

The *Complete* Universe

Russell William Ian Johnston, M.Sci. (Hons)



University
of Glasgow

Presented for the degree of
Doctor of Philosophy
The University of Glasgow
July 2009

Abstract

The work in this thesis charts the revival and development of research that tests some of the core fundamental assumptions that characterise the study of magnitude-redshift surveys for the estimation of galaxy luminosity functions (LF). Estimating LFs, either parametrically or non-parametrically, generally requires the assumption of separability between the LF, $\phi(L)$, and the density function $\rho(z)$. The work carried out initially by [Rauzy \(2001\)](#) amounts to a test statistic, T_c , constructed from the cumulative luminosity function (CLF). It is a direct probe of separability and therefore is rendered a magnitude completeness test to identify the presence of potential systematics and/or evolution within a particular survey sample.

We originally applied Rauzy's test of completeness to the Millennium Galaxy Catalogue (MGC), the Two Degree Field Galaxy Redshift Survey (2dFGRS) and the Sloan Digital Sky Survey (SDSS). We then extended the T_c statistic for data-sets characterised by two distinct faint and bright apparent magnitude limits. Following on from this we have developed a variant on T_c that we have named, T_v , which was constructed instead from the cumulative density function (CDF) and can be considered a differential form of the much celebrated, [Schmidt \(1968\)](#) V/V_{\max} statistic.

The completeness analysis of data-sets such as the 2dFGRS and the Clowes Campusano Large Quasar Group Survey (CCLQG) have also lead to developing a procedure that will optimise our estimators based on the signal-to-noise of our sampling technique.

Finally, we have developed a new, robust statistical probe to constrain evolutionary models applied to current and future redshift surveys. This probe exploits the fundamental assumption of separability coupled with a maximum entropy technique to constrain the evolutionary parameter that characterises, in particular, pure luminosity evolution.

*“These are the days of miracle and wonder,
This is the long distance call.
The way the camera follows us in slow-mo,
The way we look to us all.
The way we look to a distant constellation
That’s dying in a corner of the sky.
These are the days of miracle and wonder,
And don’t cry baby, don’t cry.”*

Paul Simon, ‘The Boy in the Bubble’, 1986

Contents

List of Tables	vi
List of Figures	xiii
Preface	xiv
Acknowledgements	xvii
List of Abbreviations	xx
1 A Selective Overview of Cosmology	1
1.1 The Road to Concordance	2
1.1.1 In the beginning...	3
1.1.2 Einstein to Friedmann-Lemaître-Robertson-Walker	4
1.1.3 Supporting evidence	7
1.2 Redshift Surveys - Past, Present and Future	9
1.2.1 Methods for measuring redshifts	9
1.2.2 The early pioneers of redshift surveys	12
1.2.3 A new chapter in surveying the Universe	13
1.2.4 Surveys... The next generation	15
1.3 The Galaxy Luminosity Function	23
1.3.1 Necessary corrections	24
1.4 Conclusions	26
2 Development of Statistical Cosmology	27
2.1 The Maximum Likelihood Estimator	28
2.2 The Traditional Non-Parametric Approaches	31
2.2.1 The <i>classical</i> approach	31
2.2.2 The V/V_{\max} test	31

2.2.3	The ϕ/Φ method	37
2.2.4	The C^- method	38
2.2.5	The stepwise maximum likelihood method	41
2.3	Emerging Methods and the Future of Estimating the LF	42
2.3.1	A semi-parametric approach	43
2.3.2	A Bayesian approach	44
2.4	Tests of Independence	46
2.4.1	Efron and Petrosian independence test	46
2.5	Conclusions	48
3	Review of the Rauzy <i>ROBUST</i> Method of Completeness	52
3.1	Constructing the R01 Completeness Statistic	52
3.1.1	Assumptions and statistical model	52
3.1.2	Defining and estimating the random variable ζ	55
3.1.3	Implementing the T_c statistic	56
3.2	Applying T_c to the Millennium Galaxy Catalogue (MGC)	57
3.2.1	The Data	58
3.2.2	Selection Limits & Cosmology	60
3.2.3	k - and evolutionary corrections	60
3.2.4	Results	61
3.3	Conclusions	64
4	Extending the Method for Doubly Truncated Surveys	65
4.1	Analysis of the 2dFGRS	65
4.1.1	The Data	65
4.1.2	Selection limits	66
4.1.3	The k - and evolutionary corrections	67
4.1.4	Initial Results	67
4.2	Generalising the T_c statistic	80
4.2.1	Re-defining the random variable ζ	80
4.2.2	Estimating ζ and computing the T_c statistic	83
4.3	Applying the Revised T_c to MGC and 2dFGRS	85
4.4	Analysis of the Sloan Digital Sky Survey - Early Types	89
4.4.1	The data	89
4.4.2	Selection limits and cosmology	91
4.4.3	Results	91

4.5	Conclusions	92
5	Introducing the T_v statistic	94
5.1	Construction the T_v statistic	94
5.1.1	Defining the random variable τ	94
5.1.2	Estimating τ and computing the T_v statistic	96
5.2	Application of the T_v Statistic	97
5.3	Conclusion	99
6	Analysis of the CCLQG Survey: GALEX Selected Sample	100
6.1	The Data and Sample Selection	100
6.2	Results	101
6.2.1	Photometric redshift truncation effect	101
6.3	Final Completeness Assessment of the Data	107
6.4	Conclusions	111
7	Optimising T_c and T_v	113
7.1	Current Issues with the ζ and τ Estimators	114
7.1.1	Shot noise - the ‘flat-line’ effect	114
7.1.2	Variation in m_{lim} effect	116
7.2	Optimising the ζ and τ Estimators	126
7.2.1	Approach #1: Fixing N_{gal}	126
7.2.2	Approach #2: Measuring the signal-to-noise	130
7.3	Conclusions	153
8	Creating Mock Galaxy Catalogues	154
8.1	Different Methodologies	155
8.1.1	Theoretician’s approach	155
8.1.2	Observer’s approach	156
8.2	Which Frame?	156
8.3	Our Mock Recipe	159
8.3.1	Sampling the magnitudes	159
8.3.2	Final selection	161
8.4	Mocking MGC, 2dFGRS and SDSS (Early Types)	163
8.4.1	MGC	163
8.4.2	2dFGRS mocks	166
8.4.3	SDSS Mocks	167

8.5	Completeness Analysis of the Mocks	172
8.5.1	Completeness of the MGC mocks	172
8.5.2	Completeness of 2dFGRS mocks	177
8.5.3	Completeness of the SDSS mocks	179
8.6	Conclusions	180
9	A New Probe for Evolution	183
9.1	Methodology to Probe Evolution	184
9.1.1	ζ and the new random variable, χ	184
9.1.2	The coefficient of correlation approach	184
9.1.3	A relative entropy approach	186
9.1.4	Implementation	187
9.2	Results Part I - Coefficient of Correlation	191
9.2.1	MGC - Mocks (R01)	191
9.2.2	MGC- Mocks (JTH)	192
9.2.3	SDSS-Mocks	200
9.3	Results Part II - Relative Entropy	208
9.3.1	MGC-Mocks (R01)	208
9.3.2	MGC-Mocks (JTH)	210
9.3.3	SDSS-Mocks (JTH)	215
9.4	Conclusions	217
10	Discussion and Future Development	220
10.1	Initial Development of the Completeness Test	220
10.1.1	Reviving the Rauzy completeness test	220
10.1.2	The 2dF survey and double truncation	221
10.1.3	From T_c to T_v	222
10.2	Future Work: Part I - Error Estimation	222
10.3	The T_v Anomaly	223
10.4	Future Work: Part II - Optimisation	224
10.5	Using Completeness to Probe Evolution	225
10.5.1	The mocks	225
10.5.2	Testing for evolution	226
10.6	Future Work: Part III - Evolution and Other Avenues	227
10.7	Concluding Remarks	228

A	Modelling m_*	229
A.1	Methods for Corrections	229
A.2	How Should We Sample the Data?	232
A.2.1	A constant m_*	232

List of Tables

1.1	Summary of Past, Present and Future Redshift Surveys.	20
4.1	2dF-Northern Plate Information.	73
4.2	2dF-Southern Plate Information.	76

List of Figures

1.1	The Hubble diagram showing the relation between the recession velocity of galaxies with distance.	4
1.2	Modern observational evidence that supports the current Λ CDM model.	8
1.3	Breakdown of our current understanding of the distribution of matter in the Universe.	9
1.4	Redshift cones of the CFA survey	12
1.5	The IRAS PSCz redshift survey	13
1.6	The 2-degree field multi-object spectrographic system.	14
1.7	The 2dF galaxy redshift survey final data release showing approximately 220000 galaxies. The 2 wedges represent the Northern and Southern stips. Image courtesy of http://msowww.anu.edu.au/2dFGRS/	15
1.8	The SDSS telescope and multifibre spectrograph.	17
1.9	The Sloan Digital Sky Survey	18
1.10	Artists impression of the Square Kilometre Array	19
1.11	Schematic illustrating the characteristic shape of the Schechter luminosity function.	24
2.1	Schematic illustrating the construction of the traditional Schmidt (1968) V/V_{\max} test.	32
2.2	Schematic illustrating the construction for the generalised V/V_{\max} and V_e/V_a statistics for overlapping survey samples.	34
2.3	Schematic illustrating the construction of the C^- method	40
2.4	Schematic illustrating the construction of the Efron and Petrosian (1992) test of independence.	47
2.5	Schematic charting the development of all the major statistical methods that estimate the galaxy LF.	51

3.1	Schematic diagram illustrating the construction of the rectangular regions S_1 and S_2 , in the original Rauzy completeness test.	53
3.2	Schematic diagram illustrating the way in which we computationally select the S_1 and S_2 regions	58
3.3	Diagram illustrating movement of m_*^f beyond the apparent magnitude limit of a given survey	59
3.4	The MGC M - Z distribution.	61
3.5	Results for the completeness statistic T_c applied to MGC	62
3.6	The random variable ζ vs distance modulus, Z	63
4.1	Survey map of the 2dFGRS showing the digitised APM plate regions.	66
4.2	The M - Z distribution of our selected sub-set of the 2dFGRS.	68
4.3	The T_c statistic for our entire 2dFGRS sample	69
4.4	The M - Z distribution of the 2dF - Northern and Southern strips.	70
4.5	The 2dFGRS T_c results for the Northern and Southern regions	71
4.6	The 2dFGRS magnitude limit masks	72
4.7	T_c curves for the 2dF-NGP plate numbers: 781 to 853.	74
4.8	T_c curves for the 2dF-NGP plate numbers: 854 to 867.	75
4.9	T_c curves for the 2dF-SGP plate numbers: 349 to 410.	77
4.10	T_c curves for the 2dF-SGP plate numbers: 411 to 473.	78
4.11	T_c curves for the 2dF-SGP plate numbers: 474 to 537.	79
4.12	Plot demonstrating how a well defined bright limit can affect the resulting T_c curve	80
4.13	In this schematic we have imposed a bright magnitude limit of $m = 16.0$ mag to the MGC data-set to illustrate how the the random variables M and Z remain unseparable under the traditional construct of S_1 and S_2	81
4.14	Diagram illustrating the construction of the rectangular regions S_1 and S_2 , defined for both faint and bright limiting magnitudes.	82
4.15	2dFGRS M - Z distribution showing our adopted apparent magnitude bright limit of 13.60 mag.	85
4.16	Performance of the T_c statistic applied to our 2dFGRS sample.	86
4.17	Our generalised version of T_c applied to 2dFGRS showing the effect of varying values of δZ	87
4.18	MGC M - Z distributions with progressively fainter bright limits	88
4.19	The T_c statistic computed for the MGC survey (without $(k+e)$ -corrections)	89

4.20	The M - Z distribution for the SDSS-Early Types.	90
4.21	Performance of the T_c statistic applied, as an illustrative example, to the SDSS- Early Type elliptical galaxies.	92
5.1	Schematic diagram illustrating the construction of the rectangular regions S_3 and S_4 , defined for a typical galaxy at (M_i, Z_i) , which feature in the estimation of our new completeness test statistic, T_v	95
5.2	Comparison of the T_c and T_v statistics computed for MGC, SDSS-Early Types and 2dFGRS	98
6.1	M - Z distribution of the CCLQGS data-set.	102
6.2	Initial T_c and T_v results of the CCLQGS data. We have assumed initially that there is no well defined bright limit and as such applied the original R01 statistic.	103
6.3	Initial T_c and T_v results of the CCLQGS data. We have assumed initially that there is no well defined bright limit and as such applied the original R01 statistic.	104
6.4	Results for the CCLQGS T_v statistic with random noise added to the raw redshift distribution.	105
6.5	M - Z distribution for CCLQGS illustrating where T_v departs from T_c in Figure 6.3.	106
6.6	Completeness results when splitting the CCLQGS data in redshift slices: $1.0 < z < 2.5$	108
6.7	Construction of T_c and T_v to illustrate the bias introduced into the T_v estimator due to truncated redshifts.	109
6.8	M - Z distribution for the CCLQGS data showing the average completeness magnitude limit derived from T_c and T_v	110
7.1	T_c and T_v results for the 2dFGRS applying the JTH generalisation for values in the range of $0.001 < \delta Z, \delta M < 0.08$	115
7.2	T_c and T_v results for SDSS-Early Types applying the JTH generalisation for values in the range of $0.001 < \delta Z, \delta M < 0.08$	117
7.3	T_c and T_v results for MGC applying the JTH generalisation for values in the range of $0.001 < \delta Z, \delta M < 0.08$	118
7.4	T_c and T_v results for the CCLQG data between $0.001 < \delta M, \delta Z < 0.08$	119
7.5	Schematic illustrating the cause of the flat-lining effect (within the 3σ confidence limits) observed for small values of δZ and δM	120

7.6	T_c and T_v plots for the MGC, SDSS-Early Types and the 2dFGRS using the R01 method to illustrate the flat-line effect.	121
7.7	T_c and T_v results for 2dFGRS applying the JTH method for values in the range of $0.1 < \delta Z, \delta M < 3.0$	122
7.8	T_c and T_v results for the SDSS-Early Types applying the JTH generalisation for values in the range of $0.1 < \delta Z, \delta M < 2.0$	123
7.9	T_c and T_v results for the CCLQG data between $0.1 < \delta M, \delta Z < 4.0$. . .	124
7.10	T_c and T_v results for the MGC applying the JTH generalisation for values in the range of $0.1 < \delta Z, \delta M < 4.0$	125
7.11	MGC T_c and T_v results applying a modified version of the R01 method where the number of galaxies, N_{gal} , in the respective $S_1 \cup S_2$ and $S_3 \cup S_4$ regions is kept constant.	127
7.12	Resulting δZ and δM distribution for the respective MGC T_c and T_v statistics.	128
7.13	CCLQG T_c and T_v results applying a modified version of the R01 method where the number of galaxies, N_{gal} , in the respective $S_1 \cup S_2$ and $S_3 \cup S_4$ regions is kept constant.	129
7.14	Schematic illustrating a procedure to calculate the signal-to-noise for ζ	132
7.15	MGC survey signal-to-noise for the R01 ζ and τ estimators.	133
7.16	CCLQG survey signal-to-noise for the R01 ζ and τ estimators.	134
7.17	3D representation of the MGC signal-to-noise map derived from the application of the JTH T_c statistic.	137
7.18	2D representation of the MGC signal-to-noise map derived from the application of the JTH T_c statistic.	138
7.19	3D representation of the MGC signal-to-noise map derived from the application of the JTH T_v statistic.	139
7.20	2D representation of the MGC signal-to-noise map derived from the application of the JTH T_v statistic.	140
7.21	3D representation of the CCLQG signal-to-noise map derived from the application of the JTH T_c statistic.	141
7.22	2D representation of the CCLQG signal-to-noise map derived from the application of the JTH T_c statistic.	142
7.23	3D representation of the CCLQG signal-to-noise map derived from the application of the JTH T_v statistic.	143
7.24	2D representation of the CCLQG signal-to-noise map derived from the application of the JTH T_v statistic.	144

7.25	3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_c statistic.	145
7.26	2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_c statistic.	146
7.27	3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.	147
7.28	2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.	148
7.29	3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.	149
7.30	2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.	150
7.31	3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.	151
7.32	2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.	152
8.1	Schematic example of how we sample the absolute magnitudes for our mock catalogue from the CDF of a Schechter function.	161
8.2	Plots tracing the k-correction models as applied to MGC, 2dFGRS and SDSS (Early Types)	164
8.3	Apparent and absolute magnitude distributions comparing our MGC mocks to the real data.	165
8.4	Apparent and absolute magnitude distribution for a trial 2dFGRS mock compared to the survey subsample.	167
8.5	Histograms of apparent and absolute magnitudes comparing our 2dFGRS mocks to the real data	168
8.6	M - Z distributions for one of our MGC, 2dFGRS and SDSS mocks compared to their respective survey samples.	170
8.7	Apparent and absolute magnitude distribution for a trial SDSS (Early Types) mock compared to the survey subsample.	171
8.8	T_c and T_v results for MGC mocks in both the G-Frame and the O-Frame.	173
8.9	T_c and T_v distributions determined from the MGC mocks, applying the R01 method.	174
8.10	T_c and T_v distributions determined from the MGC mocks, applying the JTH method for $\delta Z, \delta M = 0.01$	175

8.11	T_c and T_v distributions determined from the MGC mocks, applying the JTH method for $\delta Z, \delta M = 1.0$	176
8.12	T_c and T_v results for 2dFGRS mocks in both the G-Frame and the O-Frame applying the R01 methodology.	178
8.13	T_c and T_v results for 2dFGRS mocks in both the G-Frame and the O-Frame applying the JTH methodology.	179
8.14	T_c and T_v results for SDSS mocks in both the G-Frame and the O-Frame applying the JTH methodology.	180
9.1	Example illustrating a typical (ζ, χ) distribution for a <i>complete</i> dat-set taken at the $m_* = m_{\text{lim}} = 20.0$ mag and $m_* = 20.5$ mag of an MGC mock catalogue	185
9.2	Example of measuring the entropy of a typical (ζ, χ) distribution for a <i>complete</i> dat-set taken at the $m_* = m_{\text{lim}} = 20.0$ mag of an MGC mock catalogue	186
9.3	Flow diagram summarising the implementation of our test for evolution.	190
9.4	Trial $\hat{\beta}$ vs the correlation coefficient ρ for a series of MGC mock catalogues using the R01 method for estimating ζ	192
9.5	Plot showing how varying sizes of δZ effect the ability of the ρ estimator to constrain $\hat{\beta}$	193
9.6	MGC $\hat{\beta}$ distribution at $\rho = 0$ for 200 and 1000 bootstraps.	195
9.7	200 MGC bootstraps applying the R01 method for four values of $\beta_{\text{true}} = 0.0, 1.0, 2.0$ and 3.0	196
9.8	MGC bootstrap $\hat{\beta}$ distribution at $\rho = 0$	197
9.9	ζ - χ distribution for one of the MGC mocks.	198
9.10	M - Z distribution for the MGC mock catalogue examined in Figure 9.9 with $\beta_{\text{true}} = 0$	199
9.11	Four plots demonstrating the shot noise behaviour in the SDSS mocks for varying sizes of δZ	201
9.12	Plots showing the resulting correlation coefficient, ρ , for SDSS mock catalogues generated from an evolving LF with $\beta_{\text{true}} = 0.0, 1.0, 2.0, 3.0$	202
9.13	ρ versus $\hat{\beta}$ for 200 SDSS bootstraps with $\beta_{\text{true}} = 1.0$ and 2.0	203
9.14	$\hat{\beta}(\rho = 0)$ distribution for 200 SDSS bootstraps at $\beta_{\text{true}} = 0.0$ and 1.0	204
9.15	ρ versus $\hat{\beta}$ for 200 SDSS bootstraps with $\beta_{\text{true}} = 2.0$ and 3.0	205
9.16	$\hat{\beta}(\rho = 0)$ distribution for 200 SDSS bootstraps at $\beta_{\text{true}} = 2.0$ and 3.0	206
9.17	M - Z distribution for SDSS mock with $\beta_{\text{true}} = 0$	207

9.18	Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks for varying mesh sizes applying the R01 method.	209
9.19	Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks for three different values of β_{true}	210
9.20	Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks implementing the JTH for different values of δZ	211
9.21	Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks for three different values of β_{true}	212
9.22	$\hat{\beta}(S_{\min})$ distribution for varying increments of $\hat{\beta}$	213
9.23	$\hat{\beta}(S_{\min})$ distribution for 200 MGC bootstraps for $\beta_{true} = 0, 1.0, 2.0$ and 3.0	214
9.24	Relative entropy, S versus trial $\hat{\beta}$ using SDSS mocks for three different values of δZ . and cell sizes.	215
9.25	Relative entropy, S versus trial $\hat{\beta}$ using SDSS mocks for three different values of β_{true}	216
9.26	$\hat{\beta}(S_{\min})$ distribution for 200 SDSS bootstraps for $\beta_{true} = 0, 1.0, 2.0$ and 3.0	217
A.1	MGC survey M - Z distributions with varying methods for correcting the data.	234
A.2	MGC survey T_c and T_v statistics resulting from applying varying methods for correcting the data.	235

Preface

The role of galaxy redshift surveys is now central to modern cosmology. Since their inception just over thirty years ago, they have not only provided the cosmology community with an unprecedented amount of data, but have allowed sophisticated 3D mapping of the Universe, leading to powerful constraints in areas such as large-scale structure, galaxy formation and evolution, and the nature of dark matter and dark energy. As a result there has also been an explosion in the new field of ‘statistical cosmology’ that has paved the way for the ongoing development of sophisticated statistical processes to not only test the quality of the data but also constrain and explore, for example: current evolutionary models derived from galaxy redshift surveys, baryonic acoustic oscillations and the neutrino distribution to name a few.

In the opening chapter we give a very selective overview of cosmology, highlighting the observational evidence that supports the current favoured Λ CDM concordance model. We then trace the history of galaxy redshift surveys before focusing one of their more fundamental applications in Chapter 2 - estimating the galaxy luminosity function (LF).

Estimation of the LF of extragalactic objects presents a considerable challenge due to the complexities of the many selection effects caused by detection thresholds in apparent magnitude, colour and surface brightness, coupled with systematics and/or evolutionary effects. In order to overcome these effects numerous parametric and non-parametric statistical methods have been devised, including the [Schmidt \(1968\)](#) $1/V_{\text{max}}$ estimator, the [Lynden-Bell \(1971\)](#) C^- method and various Maximum Likelihood Estimators (e.g. [Sandage et al., 1979](#)), where a parametric form of the LF is assumed. We provide a thorough review of all these major statistical developments and discuss some of the more significant variations and extensions that have arisen due to the myriad of complexities of the many redshift catalogues that are now available to us. We also include discussions of some very recent developments in this area that propose new, and perhaps improved, methods by which the LF may be determined.

Towards the end of the second chapter we introduce an area research that has

formed the backbone of the work carried out in this thesis: tests of independence. In all the above methods for estimating the LF it is fundamentally assumed that the LF is statistically independent of the 3D spatial distribution of objects for a given survey catalogue. In particular, we discuss a paper published by [Efron and Petrosian \(1992\)](#) where they present a statistical tool that can be used to test the validity of this assumption. In Chapter 3 we summarise the work carried by Stéphane Rauzy ([Rauzy, 2001](#)) which extends the Efron and Petrosian test as a method for assessing completeness in magnitude-redshift galaxy surveys. The Rauzy method constructs a statistic based on the cumulative luminosity function (CLF), called T_c , which is used to determine the *true* apparent magnitude limit of a survey that is flux-limited. We apply Rauzy’s test to the Millennium Galaxy Catalogue (MGC) and discuss the results. In Chapter 4 we examine the Two Degree Field Galaxy Redshift Survey (2dFGRS) and a sample from the Early Types catalogue of the Sloan Digital Sky Survey (SDSS). These samples were key to extending the Rauzy method:

1. for surveys which have well defined faint *and* bright apparent magnitude limits and,
2. for surveys where the presence of a bright limit is less clearly defined and therefore may be more difficult to detect.

We extend the method further in Chapter 4 by introducing a variant on the Rauzy’s T_c statistic which we have named, T_v . Unlike the T_c statistic, T_v is constructed from the cumulative distribution function (CDF) of the redshift distribution and represents a differential version of the classic [Schmidt \(1968\)](#) V/V_{\max} test. We apply T_v to all three surveys, MGC, 2dFGRS and SDSS-Early Types and compare the results with the T_c statistic.

In Chapter 5 we apply T_c and T_v to a Galaxy Evolution Explorer (GALEX) selected sample from the Clowes-Campusano Large Quasar Group Survey (CCLQG). This dataset proved to be useful in identifying an adverse effect on the T_v statistic introduced by the highly rounded photometric redshift data. Other results from this survey and the 2dFGRS have led to our ongoing research that we detail in Chapter 6. This work explores optimising both the T_c and T_v statistics by basing their calculations on a signal-to-noise threshold which will lead to a more optimised measure of completeness that overcomes the effects of shot-noise.

In the final chapters we shift our focus from completeness to developing statistical tools to probe evolutionary models applied to current survey data. In Chapter 8

we firstly introduce our methodology for creating Monte Carlo-generated mock galaxy catalogues based on the MGC, 2dFGRS and SDSS-Early Types survey samples. By constructing mock catalogues from the observed redshift distribution of these surveys, we demonstrate how they can provide us with an effective controlled testing ground without the need for the scale of computing power demanded from dark matter N-body simulations. Furthermore, we provide a full completeness analysis of all our mocks and discuss the results and implications. In Chapter 9 we apply our MGC and SDSS mock samples to test how we can exploit the fundamental properties of our statistics - where if the wrong evolutionary correction were applied, the ‘observed’ joint sampling distribution of luminosity and redshift will *not* be separable. Hence, given only the sampling distribution of T_c and T_v under the null hypothesis of independence, we can translate the computed values of each statistic for different ‘trial’ evolutionary model parameters, $\hat{\beta}$, directly into confidence limits for this parameter. This was explored in two different ways: firstly by implementing the Pearson-Product correlation coefficient, ρ , that equals zero for the correct model parameter that is equivalent to β_{true} ; and secondly by adopting a fundamental, information-theory approach that utilises a measure of the relative entropy, S , of galaxy survey data which will then be minimised for β_{true} . We discovered that in the case of the MGC mocks the correlation coefficient appeared to detect the correct evolutionary test parameter, β_{true} within a fairly broad bootstrapped error distribution. However, when applied to the SDSS mocks, ρ indicated *two* possible values of β_{true} for every trial evolutionary model applied. In contrast it was found that the relative entropy approach was most effective for determining β_{true} . For both the MGC and SDSS mocks, the relative entropy method minimised to the correct value of β for every β_{true} model applied, with a resulting narrow error distribution compared to that of the ρ estimator.

Finally, in the concluding chapter we explore the key results and discuss the scope for future development in the area of statistical cosmology.

Acknowledgements

The road to completing this work was by no means an easy one and inevitably there have been many ups and downs along the way. The decision to return to academia was not made lightly and would probably not have happened at all if it were not for the love and support first and foremost from my fiancée, Róisín. Róisín has been a constant source of emotional strength, especially over the last few months, where it has often felt like this PhD would never end. Ok, now we can get married!

Naturally, my mum deserves a huge chunk of my appreciation. She has always been supportive of my choices in life for which I am eternally grateful. Thanks also to my brother Murray and my sister in law, Wendy, for getting me involved with our band during this time. It was a wonderful adventure trying to balance my PhD with creating the greatest music that may never be heard by the masses. And lastly, on the family side of things, thank you to my stepfather Graeme for his *Sunday Times* insights.

Of course, had it not been for a fateful meeting back in 2005 with Martin Hendry and Luìs Teodoro (who obviously have a gift for recognising talent!) getting back on the academic ladder would not have been possible. Martin's enthusiasm (especially for the acetate slides I brought along from my undergraduate degree project!), support and encouragement at that meeting, certainly gave me the confidence to embark on this new journey - and I still have those old acetate slides Martin!!

Luìs, on the other hand, barely said two words to me through his poker face during that meeting, but I had a feeling that I would end up working with him. Luckily, however, his demeanour was merely a front to test my resolve - or at least, that's what I told myself! Over the last three and a bit years, Luìs has been a seemingly unlimited source of ideas with an enthusiasm for cosmology that is extremely infectious. As a result he has acted as a mentor that has inspired much of this work and I am proud to be 'Mr Bastard #1'!

I would also like to thank Luìs's office mate, Tobia Carrozi for many useful chats. Tobia was a patient listener while I ranted endlessly about my work, which, lets face it, must have bored him to tears! Oh yeh, sorry for barging into your office every other

minute of the day to confer with Luis.

Thanks to Simon Driver and Nick Cross for useful email exchanges; Procheta and Fiona for putting up with me as an office mate when I left my mouldy plates and cups around the place; Norman for his for being my personal latex and mac consultant; Cris Sabiu for no particular reason other than he didn't acknowledge me in his Masters thesis; Craig Stark, Craig MacLachlan and Hugh Potts for entertaining me every lunchtime while I waited for my microwaved chilli.

Finally I would like give special thanks to Mr Anwar Din who was a fantastic physics and maths tutor to me whilst I was at high school, and without whom...

This thesis is my own composition except where indicated in the text.

July 23, 2009

List of Abbreviations

2dFGRS	Two Degree Field Galaxy Redshift Survey
2SLAQ	2dF-SDSS LRG and QSO Survey
2MASS	Two Micron All Sky Survey
6dFGS	Six Degree Field Galaxy Survey
APM	Automated Plate Measuring (Machine)
ASKA	Australian Path Finder Square Kilometre Array
ATNF	Australia Telescope National Facility
BAO	Baryonic Acoustic Oscillations
CCD	Charge-Coupled Device
CDM	Cold Dark Matter
CfA	Centre for Astrophysics
CMB(R)	Cosmic Microwave Background (Radiation)
CNOC	Canadian Network for Observational Cosmology
DSS	Digitized Sky Survey
ESO	European Organisation for Astronomical Research in the Southern Hemisphere
GALEX	Galaxy Evolution Explorer
GAMA	Galaxy And Mass Assembly
GF	Galaxie's Frame
GDDS	Gemini Deep Survey
GOODS	Great Observatories Origins Deep Survey
HDF	Hubble Deep Field
HST	Hubble Space Telescope
HUDF	Hubble Ultra Deep Field
KAOS	Kilo Aperture Optical Spectrograph
LBG	Lyman Break Galaxy
LCRS	Las Campanas Redshift Survey
LRG	Luminous Red Galaxy

MACHO	Massive Compact Halo Object
MGC	Millennium Galaxy Catalogue
MUNICS	MUnich Near-Infrared Survey
OF	Observer's Frame
PSC	Point Source Catalogue
PSCz	Point Source Catalogue Redshift (survey)
QSO	Quasi-Stellar Object
RSAGC	Revised Shapley-Ames Galaxy Catalogue
SDSS	Sloan Digital Sky Survey
SKA	Square Kilometre Array
SSRS	Southern Sky Redshift Survey
UKST	United Kingdom Schmidt Telescope
VIMOS	Visible Multi Object Spectrograph
VLT	Very Large Telescope
WIMPS	Weakly Interacting Massive Particles
WMAP	Wilkinson Microwave Anisotropy Probe

Chapter 1

A Selective Overview of Cosmology

“Religions die when they are proved to be true. Science is the record of dead religions.”

Oscar Wilde - ‘Phrases and Philosophies for the use of the Young’, 1894

A quick historical search on the SAO/NASA Astrophysics Data System (ADS) website finds the first use of phrase, *era of precision cosmology*, attributed to [Turner \(1998\)](#). Since then, it has been used in countless cosmology publications and, most probably, countless grant applications as well! However, it is a phrase that, when first stated, clearly resonated with the cosmology community. Emerging technologies during the 1990’s, such as the Hubble Space Telescope (HST), and the Cosmic Background Explorer (COBE) had allowed not only the mapping of the large-scale structure of Universe to high redshifts, but also, for the first time, detection and measurements of the temperature fluctuations in the cosmic microwave background radiation (CMBR).

Presently, we have high precision data obtained from Wilkinson Microwave Anisotropy Probe (WMAP) that seem to support the current Λ CDM model to a high accuracy by measuring precisely these anisotropies in the CMBR. There is also the ongoing Sloan Digital Sky Survey (SDSS) that, at the end of its second phase in 2008, had mapped 230 million celestial objects including: 930,000 galaxies, 120,000 quasars and 225,000 stars within $8,400 \text{ deg}^2$ of the sky. The coming years will usher in the the next generation of both radio and optical telescopes. With the advent of the Square Kilometre Array (SKA) and its associated pathfinder projects radio astronomy is sure to witness similar advances of instrumental precision in other regions of the electro-magnetic spectrum. Such technology will be able to survey the sky on an unprecedented scale, with the

hope that we can learn about the origin and evolution of cosmic magnetism, as well carrying out strong-field tests of gravity using pulsars and black holes, and advancing our knowledge of galaxy evolution to name just a few of its scientific aims.

In the optical range there are projects such as the Panoramic Survey Telescope And Rapid Response System (Pan-STARRS) in Hawaii that will measure astrometry and photometry of the available sky several times each month (see e.g. [Kaiser, 2002](#)). Its scientific goals include detection of exoplanets and gamma-ray bursts (GRBs), stellar evolution, galaxy clustering and providing extensive weak lensing maps (to name just a few). The system will utilise an array of four telescopes of $D \sim 1.8\text{m}$ and will have a 3 degree field of view with a CCD digital camera with 1.4 billion pixels. Its design is such that $6,000 \text{ deg}^2$ will be covered per night.

The precision with which we are now able to measure the Universe has increased dramatically and has allowed cosmology research to flourish over the last 20 years or so, taking centre stage in astronomy. Ironically, however, as we will see, these modern observations, now allude to the existence of quantities aptly named as Dark Energy and Dark Matter about which we know very little, but which are the dominating forces in the Universe as we understand it currently. This chapter will map the course that has lead us from a time where cosmology was perhaps more a theoretical pursuit without much substantial observational data, to the present time where we now have almost more data than we can handle, providing us with powerful constraints on the origins and evolution of the Universe.

1.1 The Road to Concordance

Today, what we perceive as modern cosmology can be attributed to the development of Einstein's theory of General Relativity in 1915. This new way of thinking about space, time and gravity has allowed, over the twentieth century, many other pioneering theorists to propose a new model of the Universe that is now referred to as the concordance (or standard) model for cosmology.

The Concordance Model represents the most concise model to date, combining astronomical observations with theoretical predictions to explain the origins, evolution, structure and dynamics of the Universe. In its current, simplest form, the model is often labelled as Λ Cold Dark Matter, where Λ refers to a non-zero cosmological constant.

1.1.1 In the beginning...

The origins of the concordance model are rooted in the Copernican principle - a fundamental assumption proposed by Nicolaus Copernicus in the 16th century that states we do not occupy a privileged position in the Universe. Since then this has proposition been generalised within a general relativistic framework whereby there is no privileged observer in the Universe.

Stemming from this, we can build the framework of concordance which is hinged upon the Cosmological Principle and which assumes that the Universe is both homogenous and isotropic. This implies that on large scales (typically >100 Mpc), the structure of the Universe looks the same in all directions (isotropic) and the average density of matter is approximately the same in all points in the Universe (homogenous).

In the early part of the 20th century before galaxies were understood to be galaxies (they were thought to be nebulae), Bulletin # 58 was published in the Lowell Observatory by [Slipher \(1913\)](#). This paper recorded one of the first radial velocity measurements of the nebula called, M31 in Andromeda. These velocity measurements were understood to be caused by the Doppler effect and in terms of electromagnetic radiation, are referred to as either blueshift (indicating movement towards the observer) or redshift (receding from the observer).

As for the Doppler shift, we define redshift, z , as the change in wavelength (or frequency) of the emitted light from an object, in this case galaxy's, divided by the rest wavelength (or frequency) of the light. This change in the observed wavelength is due to the galaxy's movement away from the observer and expressed as,

$$z = \frac{\lambda_o - \lambda_e}{\lambda_e} = \frac{\nu_e}{\nu_o} - 1 \quad (1.1)$$

where, λ_o and ν_o are the respective observed quantities, wavelength and frequency, and λ_e and ν_e are respectively the rest (or emitted) wavelength and frequency.

Sixteen years after these initial measurements were made, [Hubble \(1929\)](#) published what would turn out to be one of the most important discoveries that would radically advance our understanding of the Universe and our place within it. Hubble had discovered a linear relation between the distance to galaxy's and their corresponding recession velocity (or redshift) and provided the first observational evidence supporting the theory of an expanding Universe (see Figure 1.1). This relation is simply cast as,

$$H_0 = \frac{cz}{d} \quad (1.2)$$

where, H_0 is the Hubble constant describing the rate of expansion, c is the speed of light, z is the redshift and d is the distance to the object.

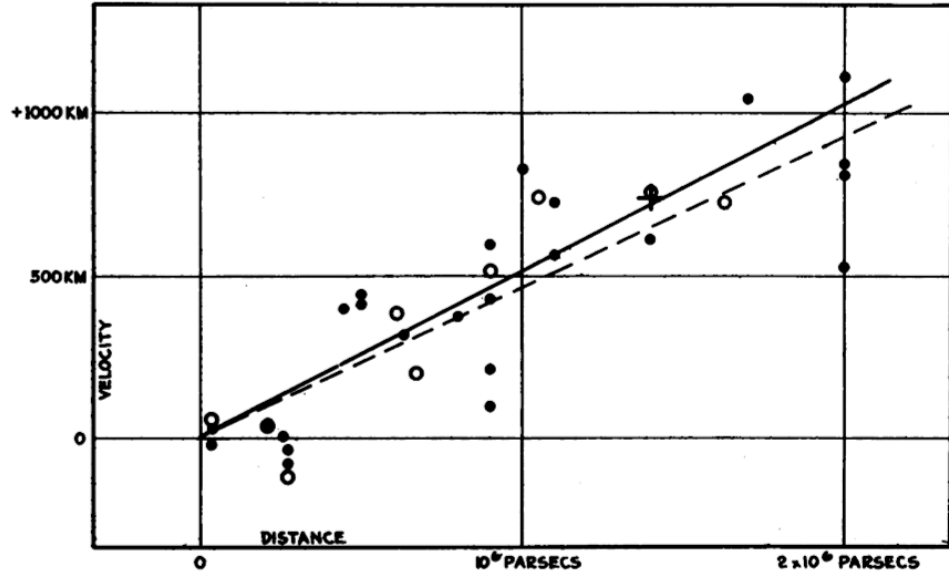


Figure 1.1: The famous Hubble diagram taken from his seminal paper, [Hubble \(1929\)](#), showing the relation between the recession velocity of galaxies with distance. This result marked the first evidence for the expansion of the Universe.

Although Hubble’s findings marked a pivotal moment in cosmology, we have to go back 14 years earlier to appreciate just how important his result was to validating, in part, the standard model.

1.1.2 Einstein to Friedmann-Lemaître-Robertson-Walker

The current standard model is formulated in the framework of Einstein’s general relativity (GR) [Einstein \(1915\)](#). Throughout this section we will attempt to refrain from carrying out a detailed derivation and simply highlight the main points. As we have already discussed, the current paradigm states that we live in a Universe based upon the cosmological principle. This principle was formally described initially by Alexander Friedman in 1922, and then extended and generalised by Howard Robertson, Arthur Walker and Georges Lemaître and is now referred to as the Friedmann-Lemaître-Robertson-Walker metric, or more simply - the FLRW metric. The metric describes geometry of space-time and it can be written as:

$$ds^2 = -dt^2 + a(t)^2 \left(\frac{dr^2}{1 - kr^2} + r^2 d\theta^2 + r^2 \sin^2 \theta d\phi^2 \right), \quad (1.3)$$

where we assume that the units for the velocity of light is $c = 1$ for the remainder of this section, $a(t)$ is the scale factor and r, θ and ϕ are co-moving spatial coordinates.

The constant, k describes the curvature of space-time and is aptly called the curvature parameter, taking on the values, -1, 0 or +1 only.

As we have said, Friedmann derived his set of equations that model the Universe from Einstein's field equation, which in their simplest form can be written as,

$$G_{\mu\nu} + \Lambda g_{\mu\nu} = 8\pi G T_{\mu\nu}, \quad (1.4)$$

where $T_{\mu\nu}$ represents the energy-momentum tensor of the matter field, $G_{\mu\nu}$ is the Einstein tensor, $g_{\mu\nu}$ is the metric tensor, Λ is the cosmological constant, G is the gravitational constant. If we assume a FLRW framework then we can describe the Universe as an ideal fluid and thus express the energy momentum tensor, $T_{\mu\nu}$ as

$$T_{\mu\nu} = (\rho + p)u_\mu u_\nu + p g_{\mu\nu} \quad (1.5)$$

where, u_μ is the 4-velocity of the matter, ρ is the mean density and p is the mean pressure. From this starting point it is possible to derive the Friedmann Equations that model the Universe based on the cosmological principle and which are given as,

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{8\pi G}{3}\bar{\rho} - \frac{k}{a^2} + \frac{\Lambda}{3} \quad (1.6)$$

and,

$$\left(\frac{\ddot{a}}{a}\right) = -\frac{4\pi G}{3}(\rho + 3p) + \frac{\Lambda}{3} \quad (1.7)$$

where, we can recast Hubble's law in terms of the scale factor yielding,

$$H = \frac{\dot{a}}{a} \quad (1.8)$$

We can now use Equations 1.6 and 1.7 to define the cosmological parameters that current observations attempt to constrain. The basic parameters can be summarised as follows:

- Ω_m - the total matter density (including dark matter)
- Ω_Λ - the dark energy density (or vacuum energy)
- Ω_k - the curvature parameter
- ρ_{crit} - the critical density
- t_0 - the age of the Universe
- w - from the equation of state

The first three parameters on the list come directly from Equation 1.6 and are dimensionless quantities such that:

$$\Omega_m = \frac{8\pi G}{3}\bar{\rho}, \quad \Omega_\Lambda = \frac{\Lambda}{3}, \quad \Omega_k = -\frac{k}{a^2} \quad (1.9)$$

Therefore, by definition,

$$\Omega_m + \Omega_\Lambda + \Omega_k = 1 \quad (1.10)$$

If we assume FLRW framework we can define the expression for the critical density, ρ_{crit} by assuming that $\Lambda = 0$ and the curvature, $k = 0$ (the critical model). From Equations 1.6 and 1.8 we get,

$$\rho_{crit} = \frac{3H^2}{8\pi G}, \quad (1.11)$$

and thus the matter density parameter, Ω_m , can be rewritten as,

$$\Omega_m \equiv \frac{\rho}{\rho_c} = \frac{8\pi G\rho}{3H^2} \quad (1.12)$$

Therefore, if we assume that our Universe is dominated by baryons, collisionless dark matter and the cosmological constant, the cosmic expansion can be described as,

$$\left(\frac{\dot{a}}{a}\right)^2 = H_0^2 \left(\frac{\Omega_m}{a} + 1 - \Omega_m - \Omega_\Lambda + \Omega_\Lambda a^2 \right). \quad (1.13)$$

As an example, in the Einstein-de Sitter model ($\Omega_m = 1, \Omega_\Lambda = 0$) the analytic solution of Equation 1.13 yields,

$$a(t) = \left(\frac{t}{t_0}\right)^{2/3}, \quad t_0 = \frac{2}{3H_0}. \quad (1.14)$$

where the present age of the Universe, t_0 , is in general given by

$$t_0 = \frac{1}{H_0} \int_0^1 \frac{r dr}{\sqrt{\Omega_m r + (1 - \Omega_m - \Omega_\Lambda)r^2 + \Omega_\Lambda r^4}}. \quad (1.15)$$

The cosmological equation of state relates to the pressure, p , and the density, ρ such that,

$$p = w\rho \quad (1.16)$$

where w is a parameter that for an accelerating Universe with a cosmological constant has the value, $w = -1$. The cosmological equation of state relates very much to the thermodynamic equation of state and the ideal gas law as it describes an isotropic universe filled with a perfect fluid.

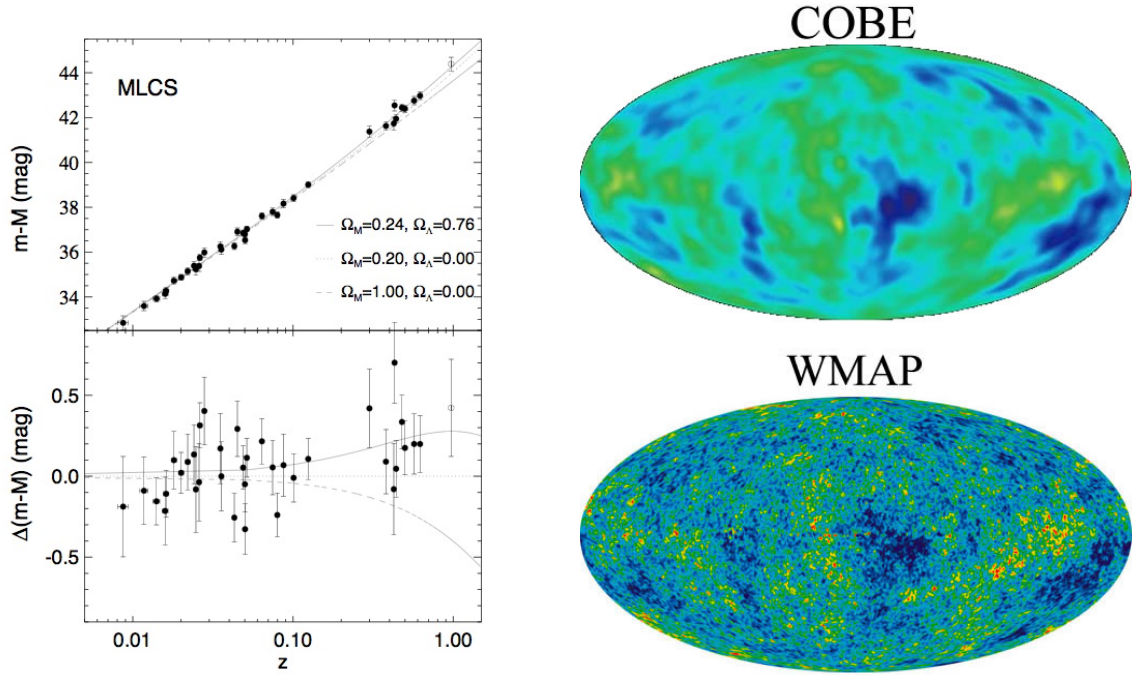


Figure 1.2: Modern observational evidence that supports the current Λ CDM model. The left panel shows the SnType Ia data from [Riess \(1998\)](#) which clearly favours a non-zero cosmological constant. The top-right image shows the temperature fluctuations recovered from the COBE satellite. The bottom-right image is represents the same fluctuations taken from the recent WMAP satellite. COBE had an angular resolution of 7° where as WMAP had an increased resolution of $\sim 0.2^\circ$. Whilst the level of detail in WMAP compared to COBE is striking, both maps show the same basic features which marked a great success for both missions. Both images courtesy of [NASA](#)

1.1.3 Supporting evidence

Apart from Hubble's discovery of an expanding Universe, there have been other key observations that support the current standard model. Probably the most significant discovery in recent times has been the strong observational evidence that supports a non-zero cosmological constant, Λ i.e. that we live in a Universe that is currently going through an accelerated expansion. Whilst it was accepted that the Universe was expanding (and even thought to have been decelerating), there was no direct evidence to support an overall cosmic acceleration. However, observations [Riess \(1998\)](#) and [Perlmutter et al. \(1999\)](#) of Type Ia supernovae provided this evidence. Type Ia supernovae are regarded to be standard candles, that is, they are considered to share the same peak absolute magnitude. Therefore measurement of the Ω_Λ parameter through the redshift-distance relation depends on comparing the apparent magnitudes of low-redshift SNe Ia with those of their high-redshift ones (see [Figure 1.2](#)).

The discovery of the Cosmic Microwave Background Radiation (CMBR) in the 1960's by Arno Penzias and Robert Wilson marked key observational evidence to support the Big Bang model. In 1992, however, another breakthrough saw results from the Cosmic Background Explorer (COBE) published by [Smoot \(1992\)](#) and [Wright \(1992\)](#) identify fluctuations in the CMBR temperature of the order of $\delta T/T \approx 10^{-5}$. Such fluctuations had already been predicted 25 years earlier (e.g. see [Sachs and Wolfe, 1967](#); [Silk, 1967](#); [Peebles and Yu, 1970](#)).

The Wilkinson Microwave Anisotropy Probe (WMAP) saw the next generation of instruments dedicated to mapping the anisotropies in the CMBR to a very high precision. However to achieve the level of detail in both maps shown in [Figure 1.2](#) required subtracting the foreground contamination from contributing sources such as our own galaxy and extra-galactic sources.

The first data-release published in [Spergel \(2003\)](#) measured the angular power spectra and supported the current standard model of a Λ -dominated universe constraining the baryon density, $\Omega_b = 0.024 \pm 0.001$, matter density, $\Omega_m = 0.14 \pm 0.02$, Hubble constant, $h = 0.72 \pm 0.05$ and $t_o = 13.7 \pm 0.2$ Gyr. The WMAP team recently released the five-year data release (see [Komatsu, 2009](#); [Dunkley, 2009](#); [Hinshaw, 2009](#); [Nolta, 2009](#); [Gold, 2009](#)) that provided, among other results, improved constraints on inflation and the curvature parameter, Ω_k ; and also provided independent corroboration of primordial nucleosynthesis or Big Bang Nucleosynthesis (BBN), which was originally explored by [Gamow \(1946\)](#).

[Figure 1.3](#) shows a simple pie chart illustrating what we know, or more appropriately, our ignorance about what makes up the Universe. Combining all these modern observations we currently estimate that dark energy is the domineering force accounting for approximately 70% of the matter of the Universe. Dark matter now plays second fiddle, accounting for 25%, whilst the baryon population is a meagre total of approximately 5%. Perhaps this chart is more of an indication of the gap between the much publicised 'era of precision' to the actual precision that we are still required to reach.

Although not a observed quantity, computer simulations have also begun to come of age. Cosmologists do not have the luxury that, say, experimental particle physicist do. They cannot re-create a Universe in the laboratory and observe the various dynamical processes as they unfold and evolve over time. However, what we can do is attempt to simulate the Universe via high performance computers. One of the milestones in computational cosmology within the last 10 years has been the development of galaxy simulations which attempt to bridge gap between the theoretical with the observed

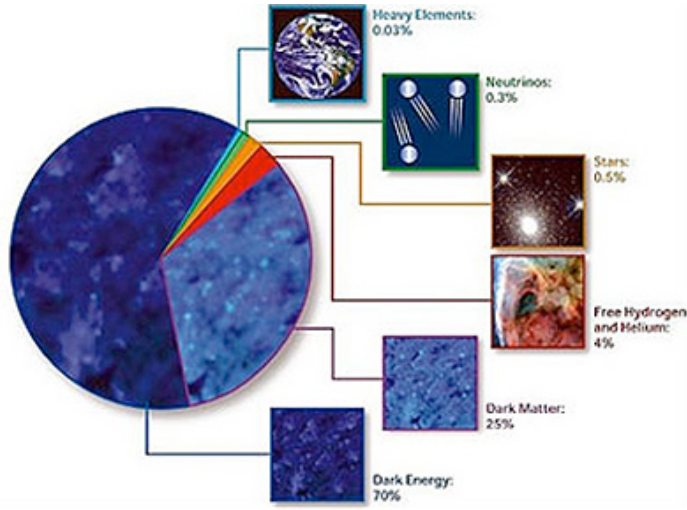


Figure 1.3: Breakdown of our current understanding of the distribution of matter in the Universe. Image courtesy of [NASA](#).

Universe. In recent years the Millennium Simulation ([Springel, 2005](#)) – a ten billion particle N-body simulation that invokes the standard Λ CDM cosmology has made it possible to compare the model predictions with the now varied and large observational surveys, which in turn allows us to probe the physical processes that drive the intrinsic evolution of galaxies.

1.2 Redshift Surveys - Past, Present and Future

Galaxy redshift surveys have played, and continue to play a vital role in our understanding of the formation, evolution and distribution of galaxies in the Universe. Prior to the 1970's, models of the structure of the Universe were based on the observed distribution of galaxies projected onto the plane of the sky. Whilst early pioneers had already identified the clustering nature of galaxies from 2-D samples (e.g. Hubble 1936 and Charlier 1922), it would take the move to three-dimensional data-sets before the wider astronomy community would accept these claims. As a consequence, this required the measurement of redshifts on a much grander scale.

1.2.1 Methods for measuring redshifts

The way in which redshifts are obtained varies from survey to survey. However, presently the most precise method of measurement is spectroscopically (z_{spec}). As we outlined earlier in this chapter, the essence of this technique involves measuring

the relative shift of a galaxy's spectral lines compared to the known position of these spectral lines as $z = 0$ i.e in the galaxy's rest frame. It is this technique that has, until very recently, been the most commonly used. Despite this, the acquisition of survey data in this way requires a two-stage process whereby one firstly measures the photometry of all galaxies in a region of space and then, secondly, targets these galaxies for spectroscopic analysis to obtain their corresponding redshift. Moreover, as telescope technologies improve we are now able to sample the deeper Universe. This has left observers reaching the limit of spectroscopic analysis since, as we move to higher redshifts, the spectral lines are redshifted out-with the spectral range. The time consuming process coupled with the physical limitation of spectroscopy has, more recently, spurred a renaissance, of sorts, in the technique of obtaining redshifts photometrically.

Photometric redshifts (z_{photo}) are obtained by identifying the shapes and broad-band spectral features of spectral energy distributions (SEDs), such as the onset of the Lyman α forest (912Å 'Lyman break'). Techniques for obtaining z_{photo} 's are by no means a recent development. In fact, during the 1960's the first method using photometry was developed for elliptical galaxies by [Baum \(1962\)](#). This technique measured the SEDs of elliptical galaxies in the Virgo cluster which were then used a calibrator and compared to the SED's from elliptical galaxies from another cluster. The displacement between two energy distributions at the 912Å Lyman-continuum discontinuity gave the redshift, and hence the distance via Hubble's Law, to the second cluster. This method represented a rudimentary version of what would become known as the 'template fitting technique' further developed by [Koo \(1985\)](#) and later by [Loh and Spillar \(1986b\)](#).

As with [Baum \(1962\)](#) the template fitting technique used the photometric data for each galaxy and converted it into the SED. By using a set of template spectra, each template was compared to the SEDs of the observed galaxies at the same 4000Å break. A χ^2 minimisation technique was then use to determine the best matching spectrum and hence the redshift of the observed galaxy. In [Loh and Spillar \(1986b\)](#) they applied their technique to 34 galaxies of known z_{spec} and in [Loh and Spillar \(1986a\)](#) they extended its use to 1000 galaxies to measure the mass density of the Universe. [Bolzonella et al. \(2000\)](#) would go on to revive this template fitting method by developing and making publically available, the *hyperz* code. Their SED fitting algorithm used a minimisation procedure based on effects such as age, metallicity, reddening and absorption in the Lyman forest.

A second, and yet, very similar technique for estimating redshifts photometrically was initially developed by [Connolly et al. \(1995\)](#), [Brunner et al. \(1997\)](#), and [Wang et al.](#)

(1998) . This approach provided an empirical technique that requires a large training set of galaxies derived from both multi-colour photometry and measured spectroscopic redshifts and is therefore, often referred to as the ‘empirical training set’ technique. A relationship is determined between the colour, C , and the redshifts, z , and, for example, in the case of Connolly et al. (1995), the redshifts are estimated via linear regression to determine a parametric fit. Whilst this remains a very popular approach, the technique is limited in redshift range by the spectroscopically obtained training sets. Moreover, both the SED fitting technique and the empirical training set techniques suffer from large photometric errors resulting in an overall accuracy of $\delta z \approx 0.1$.

The fact that galaxy surveys are now able to sample out to increasingly high redshifts, has accelerated and facilitated research into overcoming these shortcomings and improving on the large errors that plague z_{photo} . For example, work by Kodama et al. (1999) and Benítez (2000) developed a Bayesian inference approach to estimating z_{photo} where redshifts can be estimated by finding the probability $p(z|D, I)$, that is, the probability of a galaxy having a redshift z given the data, $D = \{C, m_0\}$ and the prior information, I .

Even over the last few years there seems to be accelerated interest in improving and developing new techniques. Here we mention just a few. Sheth (2007) generalises various luminosity function estimators for z_{photo} ; Ball et al. (2008) developed a nearest neighbour instance-based algorithm to improve on the z_{photo} probability density functions (PDF); Lima et al. (2008) provide an empirical method for estimating z_{photo} based on weighting the spectroscopic subsample; Stabenau et al. (2008) use galaxy surface brightness as a prior for template-based techniques; Oyaizu et al. (2008) provided a new estimator based on the training set approach using recent surveys such as SDSS and DES; Budavári (2009) present a unified approach to the template fitting and empirical techniques; Wittman (2009) use the redshift probability distribution, $p(z)$, reduce the errors; and Wolf (2009) has proposed another Bayesian approach based on the empirical training sets technique.

Although it may be unlikely that photometric redshift estimations will replace its spectroscopic counterpart completely, it seems that there are plenty of people investing the much needed time to improve on techniques to estimate them. This, at the very least, will help the continuation of photometric redshifts as a complementary option for the cosmology community, and perhaps be the replacement.

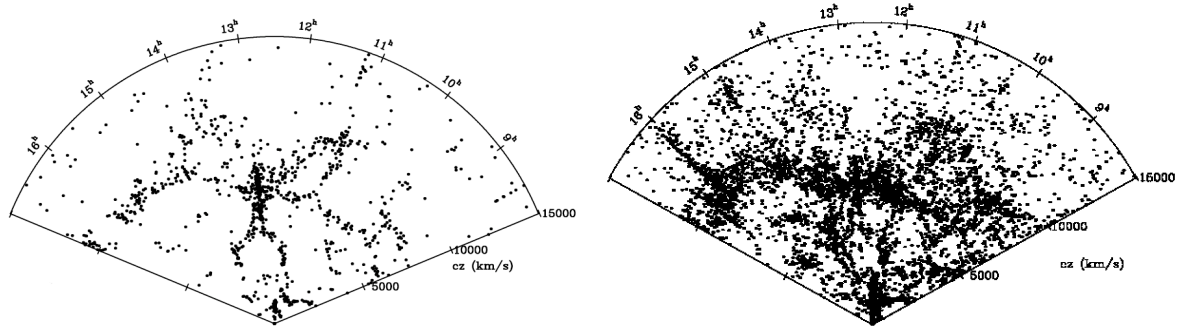


Figure 1.4: The CfA redshift surveys. The strip on the sky was 6 degrees wide and 130 degrees long with our origin being at the apex of the wedge. The initial survey, CfA, surveyed a total of 1100 galaxies as shown on the left-hand wedge. The right-hand wedge is CfA2 that surveyed a total of 18,000 galaxies in the same region as CfA. Images courtesy of the [Smithsonian Astrophysical Observatory](#)

1.2.2 The early pioneers of redshift surveys

Over the last 30 years, redshift surveys have developed into a thriving industry, and whilst they have been many and varied over the years, we will only scratch the surface, highlighting some of the more significant successes here. There is, however, a table provided with an extensive list of past, present and future (proposed) surveys in Table 1.1.

The size of redshift surveys is largely limited by the technological advances in telescopes, CCD imaging (i.e photometry) and, perhaps more importantly, spectrographs. To make a large enough survey where redshifts of thousands of galaxies could be measured would require a lot of dedicated telescope time and funding. Nevertheless, it was in 1977 that these investments were made and dedicated redshift surveys began. The first major breakthroughs in mapping large scale structure began with the CfA survey which ran from 1977 to 1982 [Huchra et al. \(1983\)](#) and measured spectroscopic redshifts for a total of 1100 galaxies out to a limiting apparent magnitude of $m_{\text{lim}} \leq 14.5$ mag (see Figure 1.4 left). This survey represented the first large area maps of large-scale structure in the nearby universe and confirmed the 3-D clustering properties of galaxies already proposed a little over 50 years previously. By extending this survey between 1985 and 1995, CfA2 ([Geller and Huchra, 1989](#)) measured a total of 18,000 redshifts out to 15,000 kms^{-1} and $m_{\text{lim}} \leq 15.5$ mag as shown in the right-hand plot of Figure 1.4. On this plot you can see more distinctly, the famous supercluster of galaxies referred to as ‘The Great Wall’. Despite this tremendous achievement, spectrographic technological constraints allowed only one galaxy at a time to be observed, making the whole

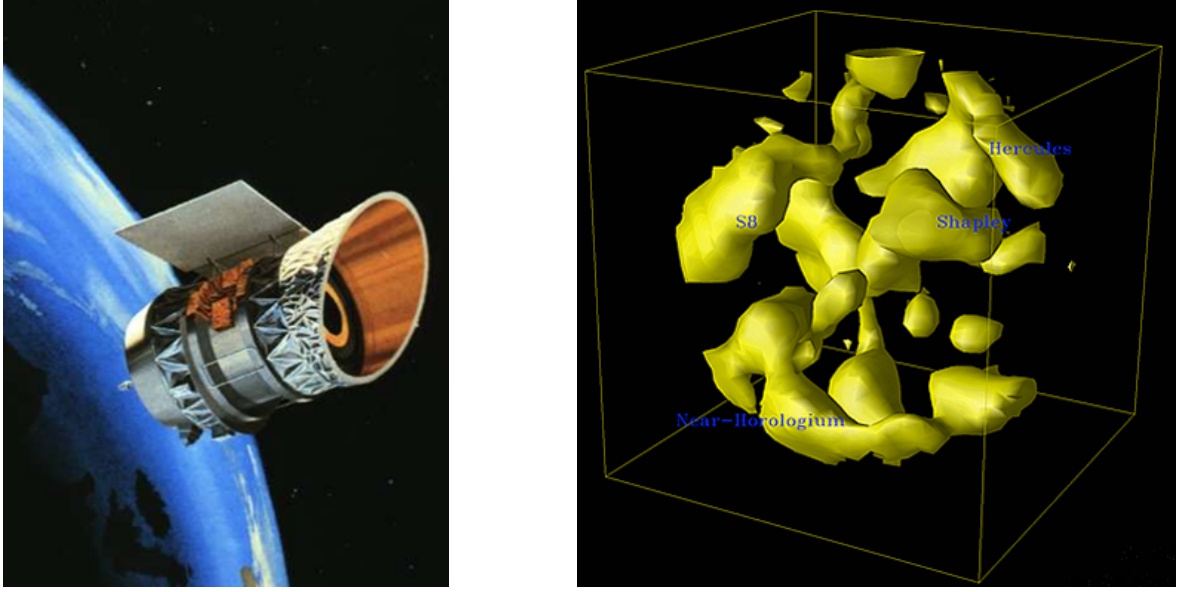


Figure 1.5: The IRAS PSCz redshift survey. The left-hand image shows the Infrared Astronomical Satellite (IRAS). The right-hand image is 3D representation of the of the PSCz survey. The left image courtesy of [NASA](#). The right image courtesy of Dr. Luis Teodoro.

process extremely time consuming. However, the 1980's saw developments in spectroscopic technology for the first multi-object fibre spectrographs that allowed between 20 and 200 galaxies to be observed simultaneously during one exposure. Moreover, the photometric technology also reached new heights.

1.2.3 A new chapter in surveying the Universe

The Infrared Astronomical Satellite (IRAS) was launched in 1983. Since effects from galactic extinction do not affect measurements in the infrared, this satellite allowed the first full-sky surveys. However, it was not until the 1990's that a new era of space-based galaxy surveys was ushered in. The PSCz redshift survey ran from 1992 to 1995 and mapped 15,411 galaxies over 84% of the sky out to 0.6 Jy [Saunders \(2000\)](#). Following this was the Hubble Deep Field-North (HDF-N) survey in 1995 ([Williams, 1996](#)). This survey utilised the Wide Field and Planetary Camera 2 (WFPC2) on the Hubble Space Telescope and for the first time, allowed unprecedented detail of faint galaxy populations to a magnitude of $m_V = 30$ mag and as deep as $z \sim 6$. In 1998 the follow up survey, HDF-South, sampled a random field in the southern hemisphere sky with equal success. Consequently a photometric redshift catalogue was created by [Fernández-Soto et al. \(1999\)](#) of 1067 galaxies in the catalogue out $z = 5.6$.

In terms of the sheer quantity of galaxies with measured redshifts the Two Degree

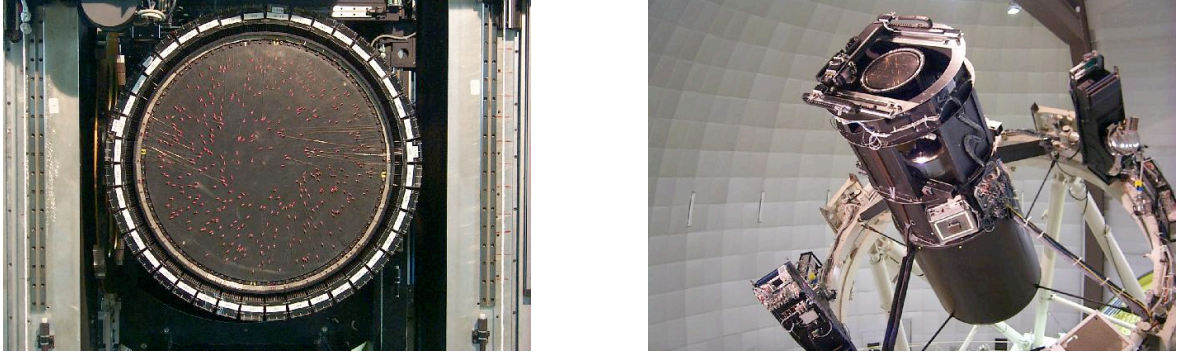


Figure 1.6: The 2-degree field multi-object spectrographic system located at the Anglo-Australian Telescope. This instrument can obtain up to 400 spectra simultaneously and marked a significant leap forward in spectrographic technology. Images courtesy of <http://www.2dfquasar.org/>

Field Galaxy Redshift Survey (2dFGRS) (Colless, 1998) could be considered the next landmark. This survey ran from 1998 to 2003 and used the multifibre spectrograph on the Anglo-Australian Telescope (AAT) which could record up to 400 galaxy redshifts simultaneously (see Figure 1.6). The photometry was taken from the APM scans of the UKST plates and measured magnitudes out to $m_{\text{lim}} = 19.45$ mag. The 2dFGRS team recovered a total of 245,591 redshifts, 220,000 of which were galaxies. With this increase in instrumentation precision and the vast number of objects catalogued, the scientific goals became equally ambitious. Some of 2dFGRS goals included measuring the power spectrum of the galaxy distribution on scales up to few hundred Mpc^{-1} , determining the galaxy LF, clustering amplitude and mean star formation rate out to a redshift $z \sim 0.5$. The survey was not only a success in terms of its achieved goals but also in the size of the collaboration: a total of 33 collaborators were involved from the UK, Australia and the USA. Figure 1.7 shows the impressive scale of the survey when compared to CfA in Figure 1.4.

At around the same time as the 2dFGRS another team was carrying out a survey called the Two Micron All Sky Survey (2MASS) (Skrutskie, 2006). This saw the return of a near-infrared full sky survey and was the first all-sky photometric survey of galaxies brighter than $m_K = 13.5$ mag and catalogued approximately 100,000 galaxies, including the ‘zone of avoidance’ - a region where dust in our own galaxy renders optical surveys near impossible.

In 2003 the HST was revisited to survey what remains today as the deepest imaging of the Universe in the optical range. The Hubble Ultra Deep Field (HUDF Beckwith, 2006) ran from 2003 to 2004 and surveyed over 10,000 galaxies out to a staggering

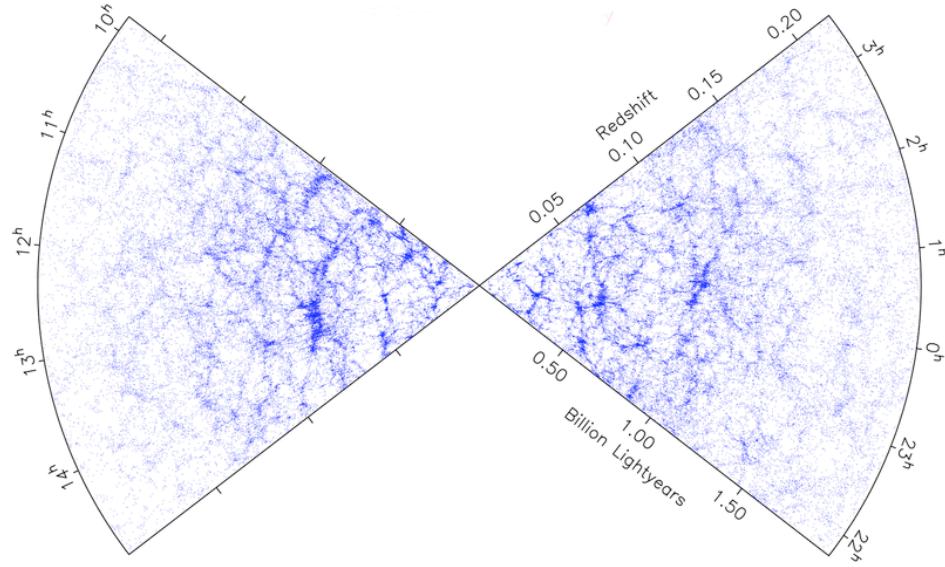


Figure 1.7: The 2dF galaxy redshift survey final data release showing approximately 220000 galaxies. The 2 wedges represent the Northern and Southern strips. Image courtesy of <http://msowwww.anu.edu.au/2dFGRS/>

$$z_{photo} = 8.$$

Moving to present day we now have the Sloan Digital Sky Survey - Data Release 7 (SDSS-DR7) which has been the most ambitious survey to date, attempting to map a quarter of the entire sky (Adelman-McCarthy, 2006). It has only just completed the second phase of the project having already catalogued just under one million galaxies. The SDSS uses a dedicated, 2.5-metre telescope on Apache Point, New Mexico, USA and a pair of spectrographs fed by optical fibres that can measure spectra and therefore the redshifts of more than 600 galaxies in a single observation.

1.2.4 Surveys... The next generation

The future of mapping the Universe looks set to begin a new ambitious phase that will truly mark the difference between 20th and 21st century observational cosmology. For example, within the foreseeable future the Square Kilometer Array (SKA) is planned to be constructed. The host country for its construction has now been shortlisted to South Africa and Western Australia. The project has been in the planning stages since 1999 (Taylor and Braun, 1999) and is one that will dwarf all previous collaborations with more than 50 institutes in 19 countries involved. SKA will operate at metre to cm wavelengths and have a field of view that will carry out instantaneous imaging of up to several tens of degrees. It will serve as an interferometric array with an overall

collecting area of order one million square metres, providing a sensitivity of around 50 times higher than the largest existing radio telescopes (see Figure 1.10). The main key science objectives include:

- Strong-field tests of gravity using pulsars and black holes;
- The origin and evolution of cosmic magnetism;
- Galaxy evolution and cosmology;
- Probing the Dark Ages;

Although the main SKA facility is due to begin construction in 2022, there are 2 pathfinder instruments due to be built in the shortlisted countries. There is the Australian SKA Pathfinder (A-SKAP) and the South African MeerKAT. Both are due to complete construction in 2012.

With this potential wealth of data cosmologists can make confident statements about the nature of the large-scale distribution of galaxies in the Universe and focus on particular galaxy types without being constrained by the low numbers which was a problem 30 years ago. For example, the VISTA Atlas survey is attempting to map $\sim 450,000$ luminous red galaxies to measure such aspects as the dark energy equation of state to a redshift $z \sim 0.7$ via baryonic oscillations.

The Large Synoptic Survey Telescope (LSST) plans to be another ambitious ground based optical telescope. With a primary mirror 10 metres in diameter, it is planned that this telescope can survey in the order of 3 billion galaxies and provide extensive weak lensing maps to probe both dark matter and dark energy. Construction is due to start in 2010 and to go online in 2015.

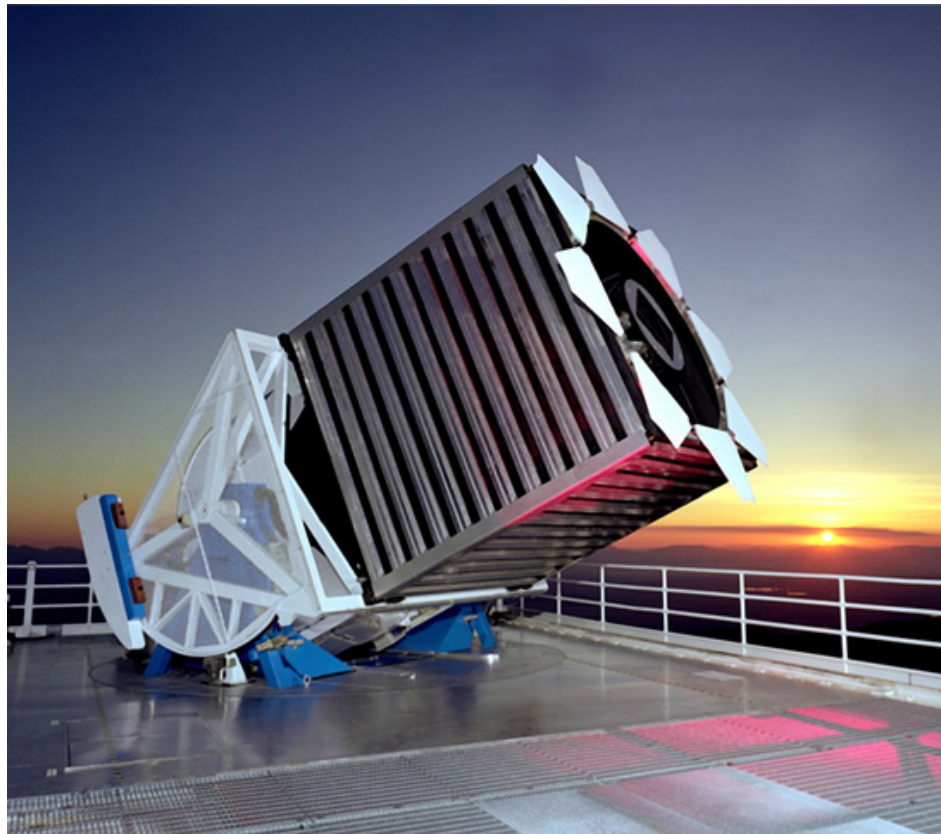


Figure 1.8: The top image shows the 2.5 m telescope located at Apache Point Observatory in New Mexico . in used for the Sloan Digital Sky Survey (SDSS). The bottom image shows the SDSS multi-object fibre spectrograph which can measure up to 600 objects simultaneously. Images courtesy of <http://skyserver.sdss.org/>

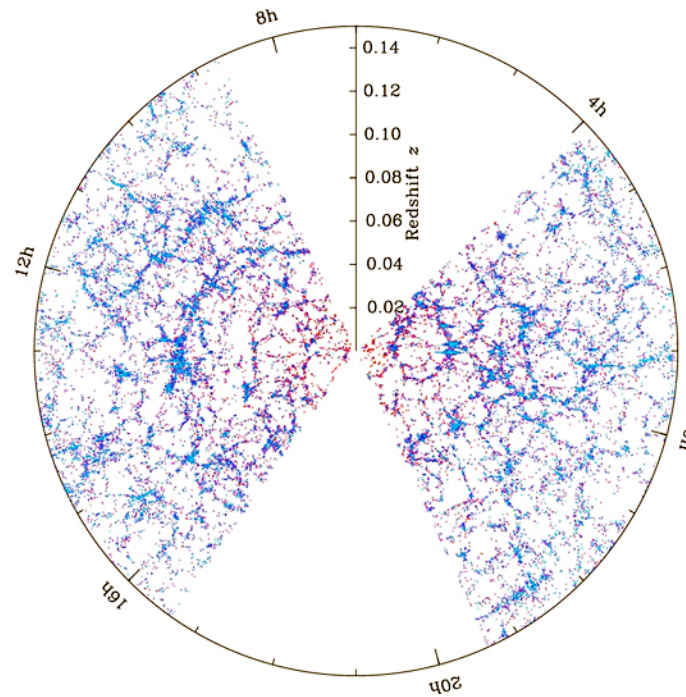


Figure 1.9: The Sloan Digital Sky Survey redshift distribution on the sky. Image courtesy of <http://www.sdss.org>.

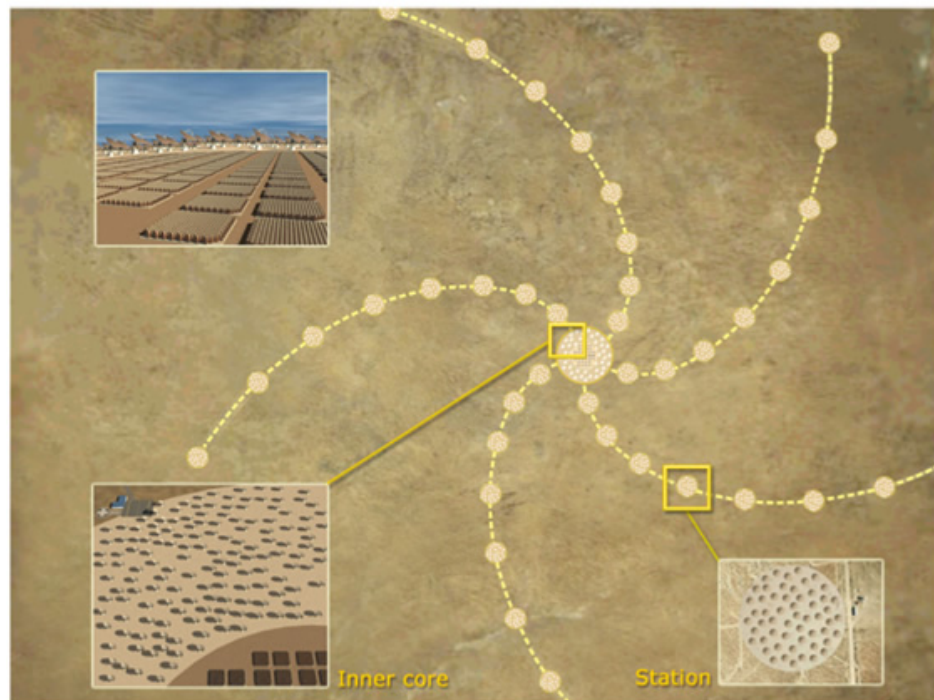


Figure 1.10: Artists impression of the Square Kilometre Array. Image courtesy of <http://www.skatelescope.org/>

Table 1.1: Summary of Past, Present and Future Redshift Surveys. See the glossary for the full names of all survey acronyms listed in the table. A brief description for smaller, perhaps lesser known surveys, has been given in place of a survey title.

Survey	No. of Objects (Galaxies/Quasars)	Coverage (deg ²)	z -limit	Source	Status
Past					
CfA-1	1,100	17,000	0.05	Huchra et al. (1983)	1977-1982
RSAGC	1191	20,649	0.05	Sandage and Tammann (1981)	1978- 1981
Survey for superclusters	~500	varied	~ 0.03	Gregory and Thompson (1978b,a, 1984)	1978-1984
Boöts void surveys	300	varied	~0.03	Kirshner et al. (1978, 1981, 1983, 1987)	1978-1987
HI 21 cm survey	2367	varied	0.0033	Fisher and Tully (1981)	1981
CfA-2	18,000	17,000	0.05	Geller and Huchra (1989)	1985-1995
CFRS	948	0.14	1.0	Lilly et al. (1995)	1995
HDF-North	3000	2.38	~ 6.0	Williams (1996)	1995
LCRS	26,418	700	~ 0.1	Shectman et al. (1996)	1988-1996
SSRS	2028	6434	0.04	da Costa (1988)	1988
Durham-UKST (DURS)	2,500	1,500	~ 0.11	Ratcliffe (1996)	1991-1996
IRAS PSCz	15,411	1,500	0.1	Saunders (2000)	1992-1999
Stromlo-APM (SAPM)	1,797	4,300	0.13	Loveday et al. (1992)	1992-2000
LBG-z3	500	0.38	~3.0	Steidel et al. (2003)	1994-2000
CNOC2	6,200	1.5	0.15 to 0.55	Carlberg (1999)	1995-1999
Eso Slice Project (ESP)	3,342	23	~0.1	Vettolani (1997)	1996-1999
2MASS	>1,000,000	37,000	0.02	Skrutskie (2006)	1997-2001
SSRS2	5,369	varied	0.04	da Costa (1998)	1998
HDF-South	2016	2.38	~ 5.6	Williams (2000)	1998
2dFGRS	220,000	1,500	~0.3	Colless (2001)	1998-2003
K20	500	0.014	1.5	Cimatti (2002)	1999-2000

Continued on next page

Table 1.1 – continued from previous page

Survey	No. of Objects	Coverage (deg ²)	z -limit	Source	Status
MGC	10,095	37.5	~ 0.46	Liske et al. (2003)	1999-2005
MUNICS	500	1.0	0.6 to 1.5	Drory (2000)	2000-2003
DEEP2	150,000	3.5	0.75 to ~ 1.0	Davis (2003)	2002-2005
GOODS	-	0.083	~ 2.0	Vanzella (2005)	2002
HUDF	10,000	0.00305	8.0	Beckwith (2006)	2003-2004
VIMOS VLT Deep Survey	150,000	16	~ 4.0	Le Fèvre (2004)	2003-2005
2SLAQ-LRG	13,000	180	~ 0.7	Cannon (2006)	2003-2006
GDDS	309	25.7	0.8 to 2.0	Abraham (2004)	2004-2007
SDSS I & II -DR1 to 7	930,000	11,663	~ 5.5 (z_{photo})	$\left\{ \begin{array}{l} \text{Abazajian (2003)DR1} \\ \text{Abazajian (2009)DR7} \end{array} \right.$	2003-2008
Present					
6dFGS	120,000	17,000	~ 0.15	Jones et al. (2005)	2001-
UKIRT	-	7500	$\left\{ \begin{array}{l} \text{ESOs: 1.0 to 2.0} \\ \text{QSOs: 7.0} \end{array} \right.$	Lawrence (2007)	2005-2012
zCOSMOS	35,000	1.0 - 1.7	$\left\{ \begin{array}{l} \text{Part I: } < 1.2 \\ \text{Part II: 1.5 to 3.0} \end{array} \right.$	Lilly (2007)	2005-2009
WiggleZ	400,000	1,000	0.5 to 1.0	Blake (2008)	2006-2010
Future					
GAMA	250,000	240	0.5	http://www.eso.org/~jliske/gama/	2009

Continued on next page

Table 1.1 – continued from previous page

Survey	No. of Objects	Coverage (deg ²)	z -limit	Source	Status
VISTA Atlas	350,000	4,500	$\begin{cases} \text{LRGs: } \sim 0.7 \\ \text{QSOs: } \sim 7.0 \end{cases}$	http://astro.dur.ac.uk/Cosmology/vstatlas/	2009
LSST	~ 3 billion	$\sim 20,000$		http://www.lsst.org/lsst	2015
SKA	1×10^9	-		Taylor and Braun (1999)	~ 2020
KAOS	~ 5 million	~ 4000		http://www.dsg.port.ac.uk/~bruce/kaos/	Proposed

1.3 The Galaxy Luminosity Function

One of the most fundamental and yet challenging problems in observational cosmology is characterising the luminosity distribution of galaxies in the Universe. Just as the study of the distribution of stellar luminosities can be a tool to help understand the physics of star formation and stellar structure, Cosmologists hope to learn about the processes of galaxy evolution by studying the distribution of galaxy luminosities. This is achieved through the determination of the optical luminosity function (LF), $\Phi(L)$, which is an essential tool for interpreting large-scale structure, determining the luminosity density and constraining galaxy formation models.

The LF describes the relative number of galaxies of different luminosities by counting the number of galaxies in a representative volume of the Universe which measures the co-moving number density of galaxies of luminosity, L , per unit luminosity such that,

$$dN = \Phi(L) \frac{dL}{L} dV \quad (1.17)$$

The quantity, $\Phi(L)$, provides us with robust handle to compare the difference between different sets of galaxies i.e. at different redshifts, galaxy types, environment etc... The LF is most commonly described by the Press-Schechter function named after its pioneers, William Press and Paul Schechter ([Press and Schechter, 1974](#)) and is typically written in the form given by,

$$\Phi(L)dL = \phi^* \left(\frac{L}{L_*} \right)^\alpha \exp \left(\frac{-L}{L_*} \right) \frac{dL}{L_*}, \quad (1.18)$$

where,

- ϕ^* is a normalisation factor defining the overall density of galaxies, usually quoted in units of, $h^3 \text{Mpc}^{-3}$.
- L_* is the characteristic luminosity which is approximately the luminosity at which most of the light of galaxies is emitted.
- α defines the faint-end slope of the luminosity function. It is typically negative, implying large numbers of galaxies with low luminosities.

Contrastingly, galaxy surveys in infra-red have yield LF's that do not seem to fit the standard Press-Schechter. For example in [Lawrence et al. \(1986\)](#) the following power law analytical form was fitted to data obtained from the Infrared Astronomical Satellite (IRAS),

$$\phi(L) = \frac{d\Phi}{dL} = \phi^* L^{1-\beta_1} \left(1 + \frac{L}{L_* \beta_2} \right)^{-\beta_2}, \quad (1.19)$$

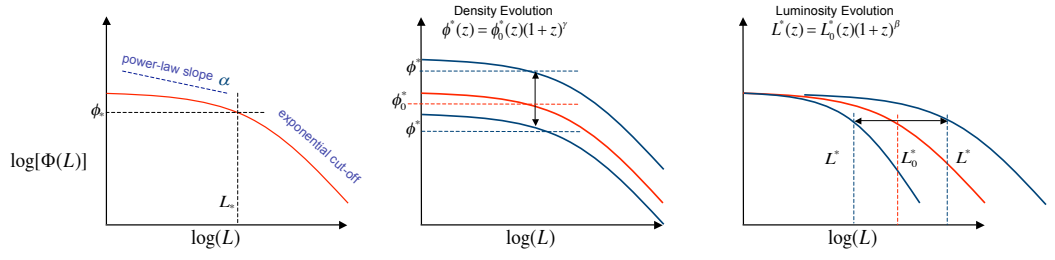


Figure 1.11: Schematic illustrating the characteristic shape of the Schechter luminosity function. The LF is typically presented in terms of log of the luminosities. The left hand plot shows the three parameter for with the Schechter LF depends upon - L_* , α and ϕ_* . The middle plot illustrates how we would expect the shape of the LF to change with the presences of pure density evolution within a survey, whilst the right-hand plot illustrates the effect of pure luminosity evolution. The degree with which both forms of evolution are effected is dependent on their respective evolution parameters γ and β .

where β_1 and β_2 define the slopes of the two power laws. In [Saunders \(1990\)](#) a log-Gaussian form was adopted for a survey also using the IRAS given by,

$$\phi(L) = \frac{d\Phi}{dL} = \phi^* \left(\frac{L}{L^*} \right)^{1-\gamma} \exp \left[-\frac{1}{2\sigma^2} \log_{10}^2 \left(1 + \frac{L}{L^*} \right) \right] \quad (1.20)$$

1.3.1 Necessary corrections

As we will go on to discuss in greater detail in the next chapter, there are several issues that affect the accuracy with which one can determine the LF. Two of the main issues relevant to this course of research are k -corrections and evolution.

1.3.1.1 The k -correction

The use of k -correction can be traced backed to early 20th century pioneering observers such as [Hubble \(1936a\)](#) and [Humason et al. \(1956\)](#), where the term, k , referred to the particular band-pass. The observed wavelength from a galaxy is different from the one that was emitted due to cosmological redshift, z . The k -correction allows us to transform from the observed wavelength, λ_o when measured through a particular filter (or bandpass) at z , into the emitted wavelength, λ_e in the rest frame at $z = 0$. Whilst there are very thorough and pedagogical reviews of its derivation e.g [Hogg et al. \(2002\)](#), we shall provide a simplified summary.

Firstly we consider the relation between the corresponding emitted frequency, ν_e

and the observed frequency, ν_o , given by,

$$\nu_e = (1 + z)\nu_o. \quad (1.21)$$

By now considering the shift of the spectrum through a particular band-pass we can relate the observed flux, $f_\lambda(\lambda_o)$ to the emitted luminosity, $L_\lambda(\lambda_e)$ by,

$$f_\nu(\nu_o)d\nu_o = \frac{L_\nu(\nu_e)d\nu_e}{4\pi d_L^2} = \frac{L_\nu(\nu_e)(1+z)}{4\pi d_L^2}d\nu_o, \quad (1.22)$$

where d_L is the cosmology dependent luminosity distance. The flux at ν_o is related to the luminosity at ν_o by,

$$f_\nu(\nu_o) = \frac{L(\nu_o)}{4\pi d_L^2} \left(\frac{\nu_e L(\nu_e)}{\nu_o L(\nu_o)} \right), \quad (1.23)$$

where, in the most general sense,

$$k(z) = \left(\frac{\nu_e L(\nu_e)}{\nu_o L(\nu_o)} \right). \quad (1.24)$$

Converting this correction from fluxes and luminosities to magnitudes yields,

$$m = M + 5 \log(d_L) + k(z), \quad (1.25)$$

where m and M are the respective apparent and absolute magnitudes of the galaxy.

1.3.1.2 Evolution

Understanding the origins and growth of structure that form the galaxies we observe today is one of the many driving forces behind current cosmological research. Whilst the intricacies of galaxy evolution are beyond the scope of our work, we do draw upon evolutionary models as applied to estimating the galaxy LF which can be described in two broad groups:

Pure Luminosity Evolution (PLE) assumes that massive galaxies were assembled and formed most of their stars at high redshift and have evolved without merging i.e. not convolved number evolution. The correction applied to account for this form of evolution is redshift and galaxy type dependent and of the form:

$$L_*(z) = L_0^*(z)(1+z)^\beta \quad (1.26)$$

converting to absolute magnitudes gives

$$M_*(z) = M_0^*(z) - 2.5\beta \log_{10}(1+z) \quad (1.27)$$

where, β is the evolution parameter and although is galaxy type dependent, a global correction is often used. The effects on the galaxy LF due to this form of evolution is illustrated in the right-hand panel of Figure 1.11.

Pure Number (Density) Evolution Number evolution assumes that galaxies were more numerous in the past but have since merged. The treatment of this type of evolution is complex involving mergers of massive halos under Cold Dark Matter (CDM) framework of N-body simulations. PDE is generally modelled as,

$$\phi_*(z) = \phi_0^*(z)(1+z)^\gamma \quad (1.28)$$

where the parameter, γ , is the number evolution parameter. As with the PLE model, we can see the effects of this form of correction on the LF in the middle panel of Figure 1.11.

1.4 Conclusions

In this chapter we have taken a very broad look at the development of modern observational cosmology. We have summarised the development of the current standard model of the Universe, Λ CDM and highlighted some of the observational measurements that support this model. We have also examined the development of redshift surveys through the ages and discussed how they have played and continue to play a key part in our continuing quest to study the origins and evolution of the Universe.

One of the main issues pertinent to the course of study throughout this thesis is the need for developing robust statistics to cope with vast amounts of data that we are now acquiring. However, we must be cautious. Despite ongoing claims that we are working in a field tagged as ‘precision cosmology’, there remains critical selection effects inherent from observations that have to be accounted for and properly understood if we are to achieve the level of required precision.

For example, redshift surveys compiled the neutral hydrogen HI 21 cm line are often used to determine the mass function in the nearby Universe. The completeness of these surveys is of major importance and difficult to quantify. Galaxy detection requires significant amounts of neutral hydrogen gas and targeting using an optical counterpart. However, low surface brightness galaxies are not easy to detect optically which can lead to significant biasing (see e.g. Zwaan et al., 1997; Salzer and Haynes, 1996; Schneider, 1996).

In the following chapter we focus in on the development of statistical methods as applied to another fundamental area of research concerning estimating the galaxy luminosity function with exploration into current tests of completeness.

Chapter 2

Development of Statistical Cosmology

*“See me in snapshots narrating my previous lives,
And a mountain of other lies.”*

‘Language of Flowers’ - Pale Saints, 1993

In this chapter we consider the development and impact of the statistical methodology within the branch of observational cosmology that attempts to determine the galaxy luminosity function (LF) derived from magnitude-redshift surveys. From the previous chapter, we will refer to the respective differential LF as $\phi(M)$, where M is the absolute magnitude of the object, and the density function $\rho(z)$, where z is the redshift of the object. As we have discussed the LF plays a vital and fundamental role in our understanding of the distribution and brightness of galaxies in large-scale structure studies and consequently is a fundamental test for galaxy formation and evolution.

However, estimating the LF for extragalactic objects can be a very tricky process where the presence of selection effects due to detection thresholds in apparent magnitude, colour and surface brightness, coupled with systematics and/or evolutionary effects, complicate the task of describing the luminosity distribution of galaxies accurately. Therefore, in order to overcome these effects and accurately determine the probability densities of $\phi(M)$ and $\rho(z)$ numerous sophisticated statistical techniques have been devised. There are non-parametric methods such as the classical number count test (e.g. [Hubble, 1936b](#); [Christensen, 1975](#)), the [Schmidt \(1968\)](#) $1/V_{\text{max}}$ test, the ϕ/Φ method (e.g. [Turner, 1979](#)), and the [Lynden-Bell \(1971\)](#) C^- method. Alternatively,

there have also been parametric methods developed (Sandage et al., 1979) based on the Maximum Likelihood Estimator (MLE) where a parametric form of the LF (most commonly that of Schechter, 1976) is assumed. There is also a non-parametric counterpart of the MLE developed by Efstathiou et al. (1988) called the Stepwise Maximum Likelihood method (SWML). Before we begin our review, we note that there are two important fundamental assumptions common to most of these methods. Firstly, they assume that the survey catalogue in question is *complete* to some specified apparent magnitude limit.

For clarity, completeness, in a cosmological context, can be defined in two distinct ways: firstly, ‘redshift’ completeness - the percentage of successfully measured redshifts over a list of targets within a survey; and secondly, magnitude completeness - the probability that a galaxy of apparent magnitude, m , is observable. Throughout this thesis we are referring to the latter.

The second assumption applies to all non-parametric and maximum likelihood estimators. The assumption is of separability between the probability densities, $\phi(M)$ and $\rho(z)$. The bivariate distribution of $P(M, z)$ where no evolution is present, can be expressed it in terms as the product of two univariate distributions such that,

$$P(M, z) = \phi(M)\rho(z) \quad (2.1)$$

There have been rigorous reviews of both parametric and non-parametric methods as they have developed over the years. One of the first reviews by Felten (1977) performed nine determinations of the LF using variations of the *classical* method. Binggeli et al. (1988) give a very comprehensive review of all non-parametric and parametric methods that had been developed up to 1988; Willmer (1997) compare the relative merits of Lynden-Bell’s C^- method with the MLE of Sandage et al. (1979) via Monte Carlo simulations; and finally, Takeuchi et al. (2000) apply the $1/V_{max}$ estimator, the C^- estimator and two variations on the MLE to simulated 2dFGRS data and Hubble Deep Field data. Therefore, the review in this chapter is very much in the same spirit, tracing the most relevant extensions and variations of the traditional approaches as well as considering more recent emerging statistical advances.

2.1 The Maximum Likelihood Estimator

We begin with the maximum likelihood estimator (MLE). As a statistical tool, the MLE is by no means a recent development. In fact it was R. A. Fisher who first pioneered the method between 1912 and 1922 (see e.g. Fisher, 1912, 1922). For a

comprehensive historical review of the MLE see [Aldrich \(1997\)](#). However, in terms of its application within the context of observational cosmology it was Sandage, Tammann & Yahil (1979), hereafter STY-MLE who pioneered this approach as applied for the determination of the LF. This is parametric technique which assumes an analytical form for the LF and therefore eliminates the binning of data. The more popular non-parametric counterpart of the MLE called the stepwise maximum likelihood (SWML) will be discussed in the following section.

If x , a continuous random variable, is described by a probability distribution function (PDF) given by,

$$f(x; \theta_1, \theta_2, \dots, \theta_k), \quad (2.2)$$

where θ represent the parameters we wish to estimate, then the likelihood function, \mathcal{L} , is given by,

$$f(x_1, x_2, \dots, x_N | \theta_1, \theta_2, \dots, \theta_k) = \mathcal{L} = \prod_{i=1}^N f(x_i; \theta_1, \theta_2, \dots, \theta_k) \quad (2.3)$$

where x_1, x_2, \dots, x_N are N number of independent observations. It is often the case that the likelihood function is expressed in terms of the logarithmic likelihood such that,

$$\Lambda = \ln \mathcal{L} = \sum_{i=1}^N \ln f(x_i; \theta_1, \theta_2, \dots, \theta_k) \quad (2.4)$$

We therefore obtain the maximum likelihood of $\theta_1, \theta_2, \dots, \theta_k$ by maximising \mathcal{L} or Λ such that,

$$\frac{\partial(\mathcal{L} \text{ or } \Lambda)}{\partial \theta_j} = 0, \quad j = 1, 2, \dots, k \quad (2.5)$$

So in the context of estimating the parameters of the LF we consider a galaxy at redshift z for which we can define the cumulative LF and thus determine the probability that the galaxy will have an absolute magnitude brighter than M as,

$$p(M|z) = \frac{\int_{-\infty}^M \phi(M') \rho(z) f(m') dM'}{\int_{-\infty}^{\infty} \phi(M') \rho(z) f(m') dM'}, \quad (2.6)$$

where $\rho(z)$ is the density function for the redshift distribution, $f(m')$ is the completeness function which for a 100% complete survey would be,

$$f(m') = \begin{cases} 1, & m_{\text{lim}}^{\text{bright}} \leq m' \leq m_{\text{lim}}^{\text{faint}} \\ 0, & \text{otherwise.} \end{cases} \quad (2.7)$$

It follows that the probability density for detected galaxies is given by the partial derivative of $P(M, z)$ with respect to M ,

$$p(M_i, z_i) = \frac{\partial P(M, z)}{\partial M} = \frac{\phi(M_i)}{\int_{M_{\text{faint}}(z_i)}^{M_{\text{bright}}(z_i)} \phi(M') dM'} \quad (2.8)$$

Note that the density functions have cancelled thus rendering the technique insensitive to density inhomogeneities. We finally maximise the likelihood to give us,

$$\mathcal{L} = \prod_{i=1}^N p(M_i, z_i). \quad (2.9)$$

Most commonly a Schechter function is assumed where the parameters that we wish to estimate are α , M^* and ϕ^* as defined in Equation 1.18 on page 23.

Marshall et al. (1983) extend the use of the MLE for quasars by simultaneously fitting evolution parameters with the luminosity function parameters. For this they test both pure density and pure luminosity models. Similarly Saunders (1990) applied the MLE method to determine the density field without knowledge of the LF. They demonstrate that by parameterising the radial density function $\rho(|\mathbf{r}|)$ they can fit it as a step function and obtain the variation on the MLE as,

$$\mathcal{L}' = \prod_{i=1}^N \frac{\rho(z_i)}{\int \rho(z_i) (dV/dz) dz}. \quad (2.10)$$

A more recent paper by Sheth (2007) revisited the STY-MLE and generalised it for the case where photometric redshifts have been used.

Although the MLE method has become more popular than other traditional non-parametric ones there are aspects not to be overlooked. As highlighted by Springel and White (1998) the MLE offers no built-in measure of goodness-of-fit. The result of which implies that nearly any functional form can be made to ‘fit’. Furthermore, the nature of the method effectively determines the slope of the LF at any point. However, if the survey sample is not complete near the apparent magnitude limit sources close to the limit will be underestimated thus making the slope of the LF underestimated (Saunders, 1990).

2.2 The Traditional Non-Parametric Approaches

2.2.1 The *classical* approach

The *classical* method, as coined by [Felten \(1977\)](#), represents the first rudimentary binned number count approach to determining the LF and was initially developed and applied by e.g. [Hubble \(1936b\)](#), [van den Bergh \(1961\)](#), [Kiang \(1961\)](#), and [Shapiro \(1971\)](#). However, as pointed out in [Binggeli et al. \(1988\)](#) the method was not formally introduced until [Christensen \(1975\)](#), [Schechter \(1976\)](#) and [Felten \(1977\)](#).

The underlying assumption of the method is that the distribution of sources within the data-set in question is spatially homogeneous. From this starting point we count the number of galaxies N within a volume V such that,

$$\Phi \equiv \frac{N}{V} \quad (2.11)$$

The volume, $V(M)$, is calculated for the maximum distance that each galaxy with an absolute magnitude, M_i , could have and still remain in the sample. As an example, [Felten \(1977\)](#) applies the following expression to calculate the volume neglecting k -corrections,

$$V(M) = \frac{4}{3}\pi \exp[0.6(m_{\text{lim}} - M_i - 25)] \times \left[E_2(0.6\alpha \ln 10) - \frac{E_2(0.6\alpha \ln 10 \csc b_{\text{min}})}{\csc b_{\text{min}}} \right] \quad (2.12)$$

where m_{lim} is the apparent magnitude limit of the survey, α and b_{min} are related to the directional-dependent galactic absorption calculation, and $E_2(x)$ is the second exponential integral ([Abramowitz and Stegun, 1964](#), Chap 5).

The number of galaxies, N , within the absolute magnitude limits of the survey, M , is binned into an arbitrary interval e.g. $(M - \frac{1}{2}\Delta M, M + \frac{1}{2}\Delta M)$ (see [Felten, 1977](#); [Binggeli et al., 1988](#)) with each bin divided by $V(M)$ to convert the histogram to units of $\text{mag}^{-1}\text{Mpc}^{-3}$ and give a binned estimation of the differential LF, $\phi(M)$.

Whilst this method is relatively straightforward to apply, its basic assumption of homogeneity is well known to be a major handicap. At the time when galaxy surveys were shallow it was common practice to exclude clustered regions such as the Virgo cluster and members of the Local Group to try and avoid biasing in the shape and thus the parameters of the LF ([Felten, 1977](#)).

2.2.2 The V/V_{max} test

A natural development from the *classical* approach is the famous V/V_{max} test which was first described by [Kafka \(1967\)](#) but more formally detailed and applied in the now

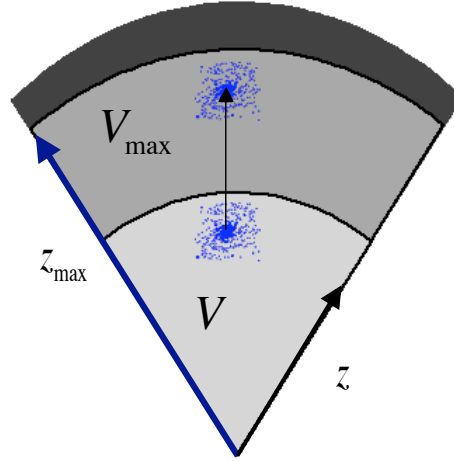


Figure 2.1: Schematic illustrating the construction of the traditional Schmidt (1968) V/V_{\max} test.

seminal paper by [Schmidt \(1968\)](#) to assess the uniformity and the cosmological evolution of quasars at high redshift (see also [Schmidt, 1972, 1976](#)). As with the *classical* method, V/V_{\max} assumes spatial homogeneity. Figure 2.1 illustrates the construction of V/V_{\max} . The basic principle of the test is simple and is defined by considering two volumes:

- V , the volume of the sphere of radius R , where R is the distance at which a galaxy was actually detected, compared to
- V_{\max} , the maximum volume within which a galaxy could have been detected. That is, the volume enclosed at maximum redshift, z_{\max} at which the galaxy in question can be observed.

It then follows that the value V/V_{\max} is expected to be uniformly distributed in the interval $[0,1]$ under the assumption of homogeneity in the survey. Thus, V/V_{\max} has expectation value

$$\left\langle \frac{V}{V_{\max}} \right\rangle = \frac{1}{2}, \quad (2.13)$$

with an often quoted statistical uncertainty of $1/(12N)^{1/2}$, where N is the total number of objects in the sample (see e.g. [Hudson and Lynden-Bell, 1991](#)). In reality, the value calculated from V/V_{\max} for a survey will deviate from $1/2$. By how much the value deviates from $1/2$ is usually considered to be either a signature of incompleteness and/or an indication of evolution: a value that is greater than $1/2$ would imply a density evolution where galaxies were more numerous in the past, whereas a value less than $1/2$ would imply that galaxies were less numerous in the past.

In the same paper, Schmidt also outlined a variation of this statistic that could be used to estimate the quasar LF,

$$\Phi = \sum_{i=1}^N \frac{1}{V_{\max}^i}. \quad (2.14)$$

Once again it was [Felten \(1976\)](#) who would dub this estimator as the ‘Schmidt’s estimator’.

2.2.2.1 Development and variations of V/V_{\max}

Since its inception, V/V_{\max} has remained a popular estimator for determining luminosity functions and as a probe of evolution, most likely due to its simplicity and ease of implementation. As a result the technique has evolved, been improved and refined over the years to accommodate the many different types of surveys that have steadily grown in size and complexity. We have selected some of the most significant developments of the method and summarised them below.

[Huchra and Sargent \(1973\)](#) were the first to extend its use to galaxies from the Markarian lists I to IV (see [Markarian, 1967, 1969a,b; Markarian et al., 1971](#)) and perform V/V_{\max} as a completeness test whilst including the Virgo Cluster and the Local Group. They showed that the effects of including such clusters had a minimal impact on the results. Furthermore, they calculated the space density $\Phi(M)$ via Schmidt’s $1/V_{\max}$ estimator, where they summed over all galaxies within an absolute magnitude interval $(M - \frac{1}{2}\Delta M, M + \frac{1}{2}\Delta M)$, where ΔM is typically 0.1.

[Felten \(1976\)](#) made an extensive comparison of $1/V_{\max}$ with the *classical* test of Equation 2.11. This paper derives a generalised version of $1/V_{\max}$ between absolute magnitude ranges $M_1 < M < M_2$ to give,

$$\int_{M_1}^{M_2} \phi dM = \sum_{i=1}^N \frac{1}{V_i} \quad (2.15)$$

and shows that it is superior to that of the *classical* estimator.

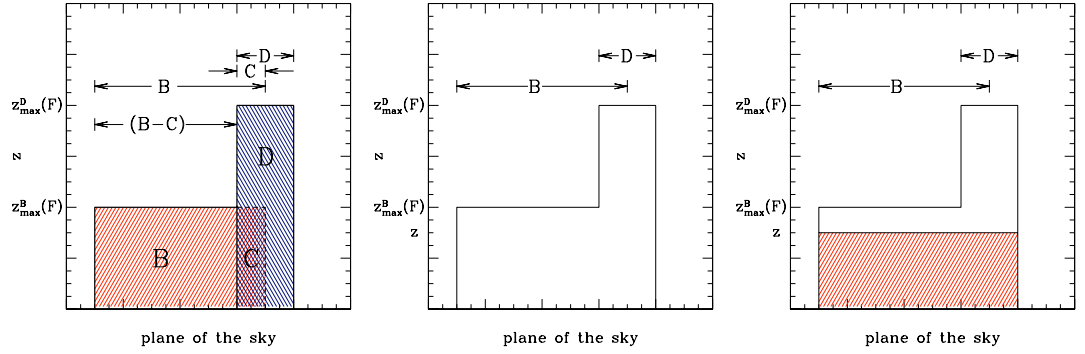


Figure 2.2: The left hand shows the construction of the generalised V/V_{\max} for two incoherent overlapping samples (B and D) into two region independent samples, (B-C) and D (see Equation 2.16). The middle panel shows the construction of V_e/V_a for a coherent sample constructed from two overlapping samples (B and D). The shaded region in the right-hand panel represents V_e for the case $z \leq z_{\max}^B(F)$ in Equation 2.18.

Avni and Bahcall (1980) generalised V/V_{\max} for multiple samples for two distinct cases:

1. Firstly, for combining independent multiple samples that are still physically separated.
2. Secondly, for combining independent samples in which the individual samples *are* physically joined together.

In the first scenario they consider complete ‘incoherent’ samples which do not overlap on the plane of the sky that could either be initially non-overlapping, or could be constructed from overlapping samples as illustrated on the left-hand panel in Figure 2.2. The term ‘incoherent’ refers to combining samples in which one remembers for each object its parent sample.

In this particular case the V/V_{\max} statistic can be constructed from overlapping samples dividing the space into tow distinct volumes (B-C) and D. For this method they show that a combined sample average of V'/V'_{\max} is given by,

$$\left\langle \frac{V'}{V'_{\max}} \right\rangle_{B-C,D} = \frac{N_{B-C}}{N_{B-C} + N_D} \left\langle \frac{V'}{V'_{\max}} \right\rangle_{B-C} + \frac{N_D}{N_{B-C} + N_D} \left\langle \frac{V'}{V'_{\max}} \right\rangle_D \quad (2.16)$$

where V' represents the density-weighted volume, N_{B-C} is the number of objects in sample (B-C) and N_D is the number of objects in sample D.

The second scenario considers the simultaneous analysis of independent complete ‘coherent’ samples in which the individual samples *are* physically joined and a new

statistic, V_e/V_a , is constructed (see illustration in the middle panel of Figure 2.2). By this description, ‘coherent’ refers to the method of combining independent samples. Here, a new variable V'_a is defined as the density-weighted volume *available* to an object for being included in the coherent sample. This new volume is defined as,

$$V'_a[F_i] = \frac{\Omega_{B-C}}{4\pi} V'[z_{\max}^B(F_i)] + \frac{\Omega_D}{4\pi} V'[z_{\max}^D(F_i)] \quad (2.17)$$

where Ω_{B-C} and Ω_D are the solid angles subtended on the sky and F_i is flux of the object. The second new variable V'_e is defined as the density-weighted volume enclosed by an object in the coherent sample and is given by,

$$V'_e[z_i, F_i] = \begin{cases} \frac{\Omega_{B-C}}{4\pi} V'(z_i) + \frac{\Omega_D}{4\pi} V'(z_i), & z_i \leq z_{\max}^B(F_i) \\ \frac{\Omega_{B-C}}{4\pi} V'[z_{\max}^B(F_i)] + \frac{\Omega_D}{4\pi} V'(z_i), & z_i > z_{\max}^B(F_i) \end{cases} \quad (2.18)$$

This first case in Equation 2.18 is illustrated in the right-hand panel of Figure 2.2. This leads to the sample average of V'_e/V'_a being defined as,

$$\left\langle \frac{V'_e}{V'_a} \right\rangle = \frac{1}{N_T} \sum_i \left\{ \frac{V'_e[z_i, F_i]}{V'_a[F_i]} \right\} \quad (2.19)$$

where N_T is the total number of objects in the two combined samples.

Hudson and Lynden-Bell (1991) recast V/V_{\max} for analysis of the diameter function of galaxies. Therefore, for diameter-limited catalogues which have both a maximum and minimum diameter cut-off they show that the completeness test can be written as,

$$\frac{V}{V_{\max}} = \frac{\theta^{-3} - \theta_{\max}^{-3}}{\theta_{\lim}^{-3} - \theta_{\max}^{-3}}, \quad (2.20)$$

where θ is the major diameter of a given galaxy, θ_{\lim} is the lower diameter limit of the survey and θ_{\max} is the maximum diameter cut-off of the survey.

Eales (1993) extended the generalised $1/V_{\max}$ estimator introduced by **Felten (1976)** to examine the evolution of the LF as a function of redshift. Similarly, **van Waerbeke et al. (1996)** looked specifically at the effects of pure luminosity evolutionary models on QSO’s via the V_{\max} estimator to constrain cosmological parameters.

Qin and Xie (1997) generalised the now familiar *Schmidt* notation in terms of a new statistic called n/n_{\max} that is applicable to *any* kind of distribution of objects.

This, therefore, would be an improved measure of the traditional V/V_{\max} test where the estimator is weighted differently and the distribution in question is assumed to homogenous. This fitting technique demonstrated that if the adopted LF is correct then the distribution of n/n_{\max} is uniform on the interval $[0,1]$,

$$\frac{n(M, z)}{n_{\max}[M, z_{\max}(M)]} = \frac{\int_0^z \Phi(M, z) dV(z)}{\int_0^{z_{\max}(M)} \Phi(M, z) dV(z)} \quad (2.21)$$

and the authors showed that its expectation value $\langle n/n_{\max} \rangle$ is $1/2$.

Following from this, another new statistic, o/o_{\max} , based on the cumulative LF and was introduced by [Qin and Xie \(1999\)](#) which is similar to but independent of n/n_{\max} .

$$\frac{o(M, z)}{o_{\max}(z)} = \frac{\int_{M_{\min}}^M \Phi(M, z) dM}{\int_{M_{\min}}^{M_{\max}(z)} \Phi(M, z) dM}, \quad (2.22)$$

This statistic is designed to provide a test of the LF and it is also shown that the two distributions of n/n_{\max} and o/o_{\max} correspond to a unique normalised LF which can provide a sufficient test for any adopted LF form. In the latter paper they apply both estimators to AAT sample data from the UVX survey [Boyle et al. \(1990\)](#).

[Page and Carrera \(2000\)](#) improve the method of constructing binned LFs using the $1/V_{\max}$ to take into account systematic errors introduced for objects close to the flux limit of a survey. As pointed out in this paper, for evolutionary studies of galaxies this traditional approach, as extended by [Avni and Bahcall \(1980\)](#) and [Eales \(1993\)](#), is very common (see e.g. [Maccacaro et al., 1991](#); [Ellis et al., 1996](#)) but can distort the apparent evolution of extragalactic populations. Through the use of Monte Carlo simulations, with a sample of 10,000 objects and simulating an unevolving two-power law model X-ray LF, they compare the $1/V_{\max}$ estimation of the differential LF given by,

$$\phi_{1/V_{\max}}(L, z) = \frac{1}{\Delta L} \sum_{i=1}^N \frac{1}{V_{\max}(i)}, \quad (2.23)$$

to their improved binned approximation of the ϕ_{est} , which assumes that ϕ does not change significantly over the luminosity and redshift intervals ΔL and Δz respectively

and is defined as,

$$\phi_{\text{est}} = \frac{N}{\int_{L_{\min}}^{L_{\max}} \int_{z_{\min}}^{z_{\max}(L)} (dV)/(dz) dz dL}, \quad (2.24)$$

where N is the number of objects within some volume-luminosity region.

[Sheth \(2007\)](#), recognising the increased use photometric redshifts in deep surveys as well as surveys vast in the number of objects, has presented an extension to the V_{\max} estimator to provide unbiased results in z_{photo} where photometric redshifts are far less precise than ones measured spectroscopically.

To do this, Sheth considers $p(z_{\text{photo}}|z_{\text{spec}})$, the probability of estimating the redshift as z_{photo} when the true value is the more accurate spectroscopic redshift, z_{spec} . He then shows the distribution, $N(z_{\text{photo}})$ of estimated redshifts is given by,

$$\frac{\mathcal{N}(z_{\text{photo}})}{dz_{\text{photo}}} = \int dz_{\text{spec}} \frac{N(z_{\text{spec}})}{dz_{\text{spec}}} p(z_{\text{photo}}|z_{\text{spec}}) \quad (2.25)$$

By considering a catalogue that has both a minimum limiting volume, V_{\min} for which an object would be too bright to be included in the catalogue, and the familiar maximum volume, V_{\max} , the number of galaxies, N , with absolute magnitude, M , in a magnitude limited catalogue is given by,

$$N(M) = \phi(M)[V_{\max}(M) - V_{\min}(M)] \quad (2.26)$$

Sheth goes on to show that the number, \mathcal{N} , of estimated absolute magnitudes, \mathcal{M} is given by,

$$\mathcal{N}(\mathcal{M}) = \int dM \phi(M) \mathcal{V}(V_{\max}, V_{\min}, M) \quad (2.27)$$

where,

$$\mathcal{V}(V_{\max}, V_{\min}, M) = \int_{D_L(M_{\min})}^{D_L(M_{\max})} dD_L \frac{dV_{\text{com}}(D_L)}{dD_L} p(M - \mathcal{M}|D_L, M) \quad (2.28)$$

where D_L is the luminosity distance and V_{com} is the comoving volume. In [Rossi and Sheth \(2008\)](#) they apply this method via mock catalogues to probe their effectiveness in areas such as the size-luminosity relation which is often distance-dependent.

2.2.3 The ϕ/Φ method

As we have previously discussed, the major drawback in the use of the V/V_{\max} is the assumption that the distribution of objects is spatially uniform. The increase in the

number and variety of redshift surveys over the years have confirmed that galaxies have strong clustering properties. Naturally, this can introduce a bias in constructing the differential LF. However, it was not long before alternative approaches were developed that could circumvent this problem.

Originally introduced by [Turner \(1979\)](#) and [Kirshner et al. \(1979\)](#) the ϕ/Φ method (as coined by [Binggeli et al., 1988](#)) is a natural progression from the *classical* method, detailed in section 2.2.1, that gives a binned estimate of the LF. For a magnitude-limited sample we calculate the ratio of the number of galaxies in the interval dM , $N(dM)$ to the total number of galaxies brighter than M , $N(\leq M)$ within the maximum volume for a complete sample.

$$\frac{N(dM)}{N(\leq M)} = \frac{dN(\leq M)}{N(\leq M)} = \frac{\phi(M)\rho(z)dM dV}{\int_{-\infty}^M \phi(M')\rho(z)dM' dV} = \frac{\phi(M)dM}{\Phi(M)} \approx d \ln \Phi(M), \quad (2.29)$$

where $\rho(z)$ is the density function and $\Phi(M)$ is the integrated LF. It is clear from this equation that the density functions cancel, thus rendering the estimator independent of the distribution of galaxies. This estimator has been further developed slightly - [Davis et al. \(1980\)](#), [Davis and Huchra \(1982\)](#) to bin the data in equal distance intervals, and more recently by [Petrosian \(2002\)](#) to account for various types of truncation. However, as shown in [de Lapparent et al. \(1989\)](#) the approximation in Equation 2.29 introduces a bias in the determination of the LF for large dM . To avoid this bias it has been common place to assume an analytical form for the LF as in [Turner \(1979\)](#), [Kirshner et al. \(1979\)](#), [Davis and Huchra \(1982\)](#) and [de Lapparent et al. \(1989\)](#). So although by its construction this estimator is non-parametric, its application has rendered it more of a parametric one.

2.2.4 The C^- method

The C^- method was introduced by [Lynden-Bell \(1971\)](#), where he applied it to the quasar data of [Schmidt \(1968\)](#). It is essentially a maximum likelihood procedure that does not require any binning of data and therefore utilises all the data in the sample. Furthermore, It has the advantage over the *classical* and $1/V_{\text{max}}$ methods as it does not require any assumptions about the distribution of objects within the data-set. This method generates a cumulative LF without normalisation with the differential LF described as a weighted sum of Dirac δ -functions.

The method can be summarised as follows. We consider the observed distribution of galaxies with absolute magnitude, M , and distance modulus, Z , plane (M, Z) , and represent the differential luminosity and density distribution functions respectively as,

$$\phi(M) = \sum_{i=1}^N \psi_i \delta(M - M_i), \quad (2.30)$$

$$\rho(Z) = \sum_{i=1}^N \rho_i \delta(Z - Z_i). \quad (2.31)$$

The distance modulus, Z , is calculated by,

$$Z = m - M = 5 \log_{10}(d_L) + 25 \quad (2.32)$$

where d_L is the luminosity distance to the object and m is the apparent magnitude. The data is then sorted such that $M_{i+1} \geq M_i$ for $i = 1, N$. We then define a region on the plane for each galaxy located at M' denoted as the $C^-(M')$ function such that,

$$C_i \equiv C^-(M_i), \quad i = 1, \dots, N \quad \left\{ \begin{array}{l} M_{\min} \leq M_i < M', \\ Z_{\min} \leq Z_i \leq m_{\lim} - M', \end{array} \right. \quad (2.33)$$

as illustrated in Figure 2.3. Therefore, the coefficients of the LF are determined from the relation,

$$\psi_{i+1} = \psi_i \frac{C_i + 1}{C_{i+1}}, \quad (2.34)$$

and can be written in the cumulative form such that:

$$\int_{M_{\min}}^M \phi(M) dM = \sum_{i: M_i < M} \psi_i \quad (2.35)$$

$$= \psi_1 \prod_{i: M_i < M} \frac{C_i + 1}{C_i} \quad (2.36)$$

Although Jackson (1974) extended the original method to account for the combining of multiple data-sets, the method remained limited to deriving only the shape of the probability density function. As a result Choloniewski (1987) revisited and improved the method to not only simplify it, but to properly normalise the LF and estimate the density evolution of galaxies simultaneously. A lesser known paper by Schmitt (1990) extended the method for samples with multi-flux limits.

Another drawback of C^- lay in its basic construction - weighted step functions - and therefore it has limited use towards the faint magnitude limit of the survey in question.

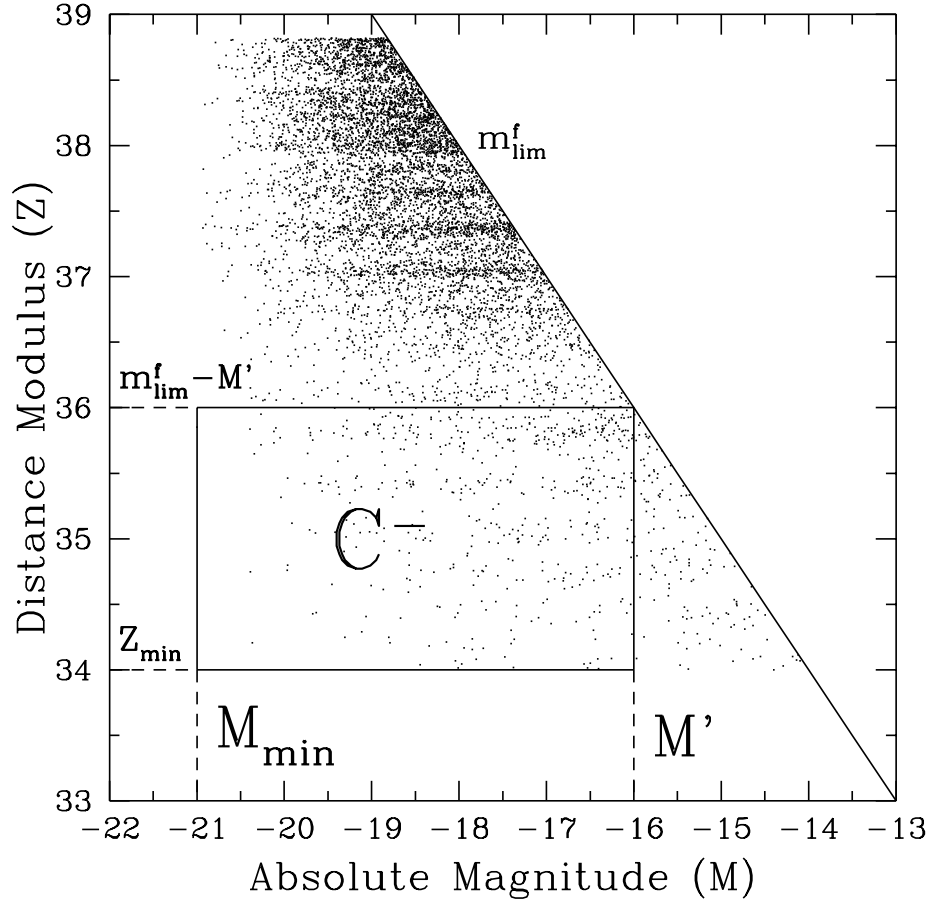


Figure 2.3: Schematic illustrating the construction of the C^- method introduced by Lynden-Bell (1971).

To overcome this, Caditz and Petrosian (1993) introduced a smoothing non-parametric method based on a Gaussian kernel, which replaces the δ -function in Equation 2.31 with,

$$K(x) = \frac{1}{(2\pi)^d |C|^{1/2}} \exp \left\{ -\frac{1}{2} \mathbf{x}^T C x \right\}, \quad (2.37)$$

where $K(x)$ is the kernel function, d is the number of measured parameters for each object, x describes the observed data and C is the inverse of the sample covariance matrix. The kernel therefore replaces the δ -function distributions given in Equations 2.30 & 2.31.

Subbarao et al. (1996) extended the method for photometric redshifts by considering, for each galaxy, the probability distribution in absolute magnitude M resultant from the photometric redshift error. By adopting a Gaussian error distribution for the function $z^*(m_i, M)$, the redshift for the i^{th} galaxy with an apparent magnitude m_i ,

they showed that for a complete magnitude limited sample the defined region $C(M)$ is now given as

$$C(M) = 0.5 \sum_i \left[\operatorname{erfc} \left(\frac{z^*(m_i, M) - z_i}{\sigma_i} \right) - \operatorname{erfc} \left(\frac{z^*(m_{\text{lim}}, M) - z_i}{\sigma_i} \right) \right], \quad (2.38)$$

where $\operatorname{erfc}(x)$ is the complementary error function.

2.2.5 The stepwise maximum likelihood method

One of the first incarnations of the stepwise MLE approach, was presented by [Nicoll and Segal \(1983\)](#) and can be thought of as a binned version of the Lynden-Bell's C^- method. This was applied in their analysis of chronometrical cosmology and also considers variations of progressive truncation in apparent magnitude as well as multivariate complete samples.

A more advanced version of Nicoll and Segal's approach by [Choloniewski \(1986\)](#) applied the same Poisson probability distribution in the MLE of [Marshall et al. \(1983\)](#) (see section 2.1). In this paper, the data is projected on the absolute magnitude M and distance modulus Z plane and divided into equal sized cells. Using the notation from [Takeuchi et al. \(2000\)](#) it is shown that the likelihood function can be represented as,

$$\mathcal{L} = \prod_{(M_i, z_j) \in S} \frac{e^{-\lambda_{ij}} \lambda_{ij}^{n_{ij}}}{n_{ij}!}, \quad (2.39)$$

where,

$$\lambda_{ij} = \frac{1}{\bar{n}} \Phi(M) \rho(x, y, z) dM dx dy dz, \quad (2.40)$$

where n_{ij} is the number of galaxies in the (i, j) cell, \bar{n} is the mean density of the sample, $\Phi(M)$ is LF, and $\rho(x, y, z)$ is the number density where the assumption of separability between M and x, y, z is employed.

The method which is actually named the stepwise maximum likelihood method (SWML) was introduced by [Efstathiou et al. \(1988\)](#) and represents the non-parametric version of the STY-MLE method. Unlike the MLE this technique does not depend on an analytical form for $\phi(M)$. Instead the LF is in effect parameterised as a series of step functions such that,

$$\phi(M) = \sum_{i=1}^N \phi_i W(M_i - M), \quad (2.41)$$

where,

$$W(x) \equiv \begin{cases} 1, & -\frac{\Delta M}{2} \leq x \leq \frac{\Delta M}{2}, \\ 0, & \text{otherwise.} \end{cases} \quad (2.42)$$

Therefore, it can be shown that the expression for the step wise likelihood is given by,

$$\mathcal{L}(\{\phi_i\}_{i=1,\dots,I} | \{M_k\}_{k=1,\dots,K}) = \prod_{k=1}^{N_{\text{obs}}} \frac{\sum_{l=1}^K W(M_l - M_k) \phi_l}{\sum_{l=1}^K \phi_l H(M_{\text{lim}}(z_k) - M_l) \Delta M} \quad (2.43)$$

where,

$$H(x) = \begin{cases} 0, & x \leq -\Delta M/2, \\ (x/\Delta M + 1/2), & -\Delta M/2 \leq x \leq \Delta M/2, \\ 1, & x \geq \Delta M/2. \end{cases} \quad (2.44)$$

The SWML has now become a very popular method for determining the LF and [Sodre and Lahav \(1993\)](#) extended it (and the STY-MLE estimator) to a bivariate distribution of magnitudes and galaxy diameters (see also [Santiago et al., 1996](#)). [Driver et al. \(2005\)](#) further extended this approach for the joint distribution of surface brightness and space density.

[Heyl et al. \(1997\)](#) extended the use of the SWML method by generalising it in a similar way as [Avni and Bahcall \(1980\)](#) did for V/V_{max} by combining various surveys with different magnitude limits, coherently. Moreover, this extension also provided an absolute normalisation and was used to probe density evolution in the LF by spectral type. [Springel and White \(1998\)](#) also explored evolution and provided a variation of the method by parameterising in terms of piecewise power laws as opposed to top-hats as in Equation 2.44. For the SDSS- Early Types [Bernardi \(2003a\)](#), replace Equation 2.42 with Gaussian weights.

2.3 Emerging Methods and the Future of Estimating the LF

Despite the continuing popularity of such methods as V/V_{max} and the SWML, there exists a need (even if it is not widely recognised at present in the wider observational cosmology community) to push the development of robust statistical techniques to the next level, so that one may overcome all the potential hazards and pitfalls inherent with current approaches already discussed in the previous sections. Therefore, we examine here two new approaches to estimating the LF that could potentially be key to satisfying such a need.

2.3.1 A semi-parametric approach

The first one by [Schafer \(2007\)](#) is by his own account statistically rigorous and considers data-sets that are truncated. There are four key points that lay the foundation for this method which are:

1. No strict parametric form is assumed for the bivariate density.
2. No assumption of independence between redshift and absolute magnitude is made.
3. No binning of data is required.
4. A varying selection function can be incorporated.

By not assuming separability Schafer decomposes the bivariate density $\phi(z, M)$ into,

$$\log \phi(z, M, \theta) = \mathbf{f}(z) + \mathbf{g}(M) + \mathbf{h}(z, M, \theta), \quad (2.45)$$

where $\mathbf{h}(z, M, \theta)$ has an assumed parametric form that folds in, for example, evolution of the LF. The functions $\mathbf{f}(z)$ and $\mathbf{g}(M)$ are determined non-parametrically. He then incorporates an extended form of the maximum likelihood approach called the ‘local’ likelihood estimator for the density estimation and applies this to 15,057 quasars from [Richards \(2006\)](#). This semi-parametric approach has the advantage of allowing the user to estimate evolution of the LF with redshift without assuming a strict parametric form for the bivariate density. The only parametric form required is that which models the dependence between redshift and absolute magnitude. We now summarise this procedure in more detail.

The local likelihood density estimation: This approach is a non-parametric extension of the MLE where one assumes the data $\mathbf{X} = (X_1, X_2, \dots, X_n)$ are observations of independent, identically distributed random variables from a distribution with density f_0 . The MLE for f_0 is defined as the $f \in \mathcal{F}$, where \mathcal{F} denotes the class of candidates for f_0 , and is maximised as,

$$\sum_{j=1}^n \log f(X_j) - n \int f(x) dx \quad (2.46)$$

From this, one can localise the likelihood criterion and thus construct the final local likelihood \hat{f}_{LL} estimator by smoothing the local estimates giving,

$$\hat{f}_{LL}(x) \equiv \left[\sum_{u \in \mathcal{G}} K^*(x, u, \lambda) \hat{f}_u(x) \right] / \left[\sum_{u \in \mathcal{G}} K^*(x, u, \lambda) \right] \quad (2.47)$$

where \mathcal{G} forms a grid $u \in \mathcal{G}$ of equally spaced values (between -3 and 3 in the authors example) of a Gaussian density with mean zero and variance of unity. The term, $K^*(x, u, \lambda)$, is therefore a kernel function such that,

$$\sum_{u \in G} K^*(x, u, \lambda) = 1 \quad \forall x.$$

By making \mathcal{G} sufficiently large, the amount of smoothing is completely dominated by the kernel function parameter, λ .

Extending to flux-limited data: This approach is extended for the use flux-limited survey data where one can include the dependence between the redshift, z , and absolute magnitude, M . A first order approximation of \mathbf{h} is made from Equation 2.45 such that,

$$\mathbf{h}(z, M, \theta) = \theta z M. \quad (2.48)$$

After an extensive derivation, a global criterion is found to be given by,

$$\begin{aligned} L^*(\mathbf{f}, \mathbf{g}, \mathbf{z}, \mathbf{M}, \theta) \equiv & \sum_{j=1}^n w_j \left(\sum_{u \in G} K^*(z_j, u, \lambda) a_u(z_j) \right. \\ & + \sum_{u \in G} K^*(M_j, v, \lambda) \mathbf{b}_v(M_j) + \mathbf{h}(z_j, M_j, \theta) \\ & - \int_{\mathcal{A}} \left\{ \exp(\mathbf{h}(z, M, \theta)) \left[\sum_{u \in G} K^*(M, v, \lambda) \exp(\mathbf{b}_v(M)) \right] \right. \\ & \left. \left. \times \left[\sum_{u \in G} K^*(M, v, \lambda) \exp(\mathbf{a}_v(M)) \right] dM dz \right\} \right), \end{aligned} \quad (2.49)$$

where $\mathbf{a}_u(z)$ and $\mathbf{b}_v(M)$ are degree p polynomials which form part of the smoothing term of K^* for local estimates. \mathcal{A} defines the region outside of which the data are truncated on the (z, M) plane.

Estimating the LF in this way has the advantage of allowing the user to estimate the evolution of the LF without assuming a strict parametric form of the bivariate density.

2.3.2 A Bayesian approach

The second technique by [Kelly et al. \(2008\)](#) adopts a Bayesian approach to estimating LFs. In this paper they derive a likelihood function of the LF that relates observed data to the true LF (assuming some parametric form). They then use a Bayesian framework to estimate the LF and the posterior probability distribution of the LF parameters via a mixture of Gaussian functions. By modelling the LF using Gaussian functions, they circumvent the problem of having to assume a parametric form.

Estimating the LF likelihood: The form of the likelihood function that they adopt for the LF estimation is derived from a binomial distribution. Whilst they highlight that the traditional approach of using a Poisson distribution is incorrect, they show that as long as the survey's detection probability is small, both approaches yield the same results. We recall the relation of the LF to the probability density of (L, z) can be written in the form of,

$$p(L, z) = \frac{1}{N} \phi(L, z) \frac{dV}{dz}, \quad (2.50)$$

where L is the luminosity, z is redshift and N is the total number of objects in the observable Universe. From this starting point, the authors assume a parametric form for $\phi(L, z)$, with parameters θ and show that the observed data likelihood function is given by,

$$p(L_{obs}, z_{obs}, \mathbf{I}|\theta, N) \propto C_n^N [p(I = 0|\theta)]^{N-n} \prod_{i \in \mathcal{A}_{obs}} p(L_i, z_i|\theta), \quad (2.51)$$

where,

$$p(L, z|\theta) = \prod_{i=1}^N p(L_i, z_i|\theta) \quad (2.52)$$

is the likelihood function for all N sources in the universe. By adding in sample selection, the probability that the survey misses a source, given by the parameters θ , is,

$$p(I = 0|\theta) = \iint p(I = 0|L, z) p(L, z|\theta) dL dz \quad (2.53)$$

In Equation 2.51 \mathbf{I} is a vector of size N taking on values:

$$I_i \begin{cases} 1 & \text{if } i^{th} \text{ source is included in survey} \\ 0 & \text{otherwise} \end{cases} \quad (2.54)$$

Finally the term C_n^N is the binomial coefficient and \mathcal{A}_{obs} is the set of n included sources.

Bayesian estimation the LF: A mixture of Gaussian functions are then used to model the LF for the joint distribution of $\log L$ and $\log z$. By taking the log of L and z fewer number of Gaussians are expected to be required to sufficiently describe the LF. Equation 2.50 is now re-written as,

$$p(L, z) = \frac{p(\log L, \log z)}{Lz(\ln 10)^2}, \quad (2.55)$$

Using a mixture of K gaussian functions, the model for the LF is given by,

$$\phi(L, z|\theta, N) = \frac{N}{Lz(\ln 10)^2} \left(\frac{dV}{dz} \right)^{-1} \sum_{k=1}^K \frac{\pi_k}{2\pi |\sum_k|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu_k)^T \sum_k^{-1} (\mathbf{x} - \mu_k) \right] \quad (2.56)$$

where $\sum_k^{-1} \pi_k = 1$, $\mathbf{x}_i = (\log L_i, \log z_i)$, μ_k is the 2-element mean position vector for the k^{th} Gaussian, \sum_k is the 2×2 covariance matrix for the k^{th} Gaussian, and \mathbf{x}^T denotes the transpose of \mathbf{x} .

The Metropolis-Hastings algorithm (MHA) is then used for obtaining random draws of the LF from the posterior distribution. Given a suitably large enough number of Gaussian functions it is flexible enough to give an accurate estimate of any smooth and continuous LF and therefore has potentially the advantage of being able to extrapolate beyond the survey flux limits.

2.4 Tests of Independence

So far we have looked at the various non-parametric and parametric methods used to determine LFs. However, as we discussed at the beginning of this chapter, the fundamental assumption for all the non-parametric techniques is the separability of the density function and the LF that allows us to write the probability distribution of observables as,

$$\Phi(L, z) = \phi(L)\rho(z) \quad (2.57)$$

This is a fundamental and crucial assumption that is generally accepted over small redshift bins (or shallow redshift surveys). However, what is sometimes overlooked is whether the assumption of separability is valid i.e. are luminosity, L , and redshift, z , of Equation 2.57 uncorrelated?

2.4.1 Efron and Petrosian independence test

In order to test this assumption of independence (or separability) between two variables, Efron and Petrosian (1992) developed a simple ranked-based non-parametric

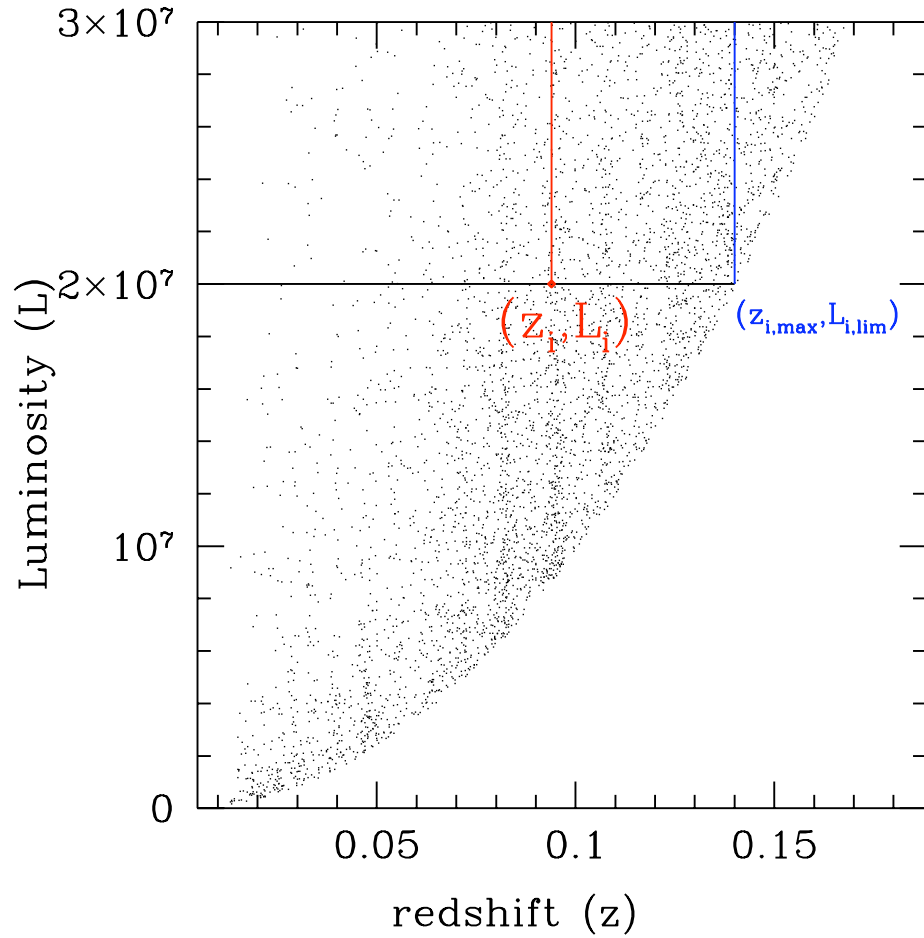


Figure 2.4: Schematic illustrating the construction of the Efron and Petrosian (1992) test of independence

test statistic for data-sets with a single truncation. Their method is, in principle, an extension of the C^- method. It is this approach which has formed the backbone of the work within this thesis. For this technique Efron and Petrosian construct the test statistic, τ , based on the rejection of independence between the random variables x and y under test. Although they give a very detailed and formal explanation, we only summarise the main points of the method here. If x and y are independent, then the rank, R_i , of x_i , for the simple case where there is no truncation, will be uniformly distributed between 1 and N with respected expectation value and variance,

$$E = \frac{1}{2}(N + 1), \quad V = \frac{1}{12}(N^2 - 1). \quad (2.58)$$

The statistic T is the defined as,

$$T_i = \frac{(R_i - E)}{V}, \quad (2.59)$$

such that R_i is now normalised to have mean of zero and a variance of unity. The hypothesis of independence is then rejected or accepted depending on the value of T_i . Moreover, one can construct confidence levels of rejection and combine T_i into a single statistic τ such that,

$$\tau = \frac{\sum_i (R_i - E)}{\sqrt{\sum_i V}} \quad (2.60)$$

where, by definition, a τ of 1 indicates a $1\text{-}\sigma$ correlation and conversely, a τ of 0 indicates that the variables are completely uncorrelated.

For the case of an imposed faint apparent magnitude limit i.e. a singly truncated data-set, the method is modified to measuring the ranks R_i for subsets of the entire set of observables. Therefore, within each set, all objects which could have been observed up to the flux limit are counted such that each set is defined as,

$$J_i \equiv j : L_j > L_i, L_{\text{lim},j} < L_i. \quad (2.61)$$

Therefore, for each galaxy an area is defined within the range $L_i < L < \infty$ and $0 < z < z_{\text{max}}(L_i)$ as illustrated in Figure 2.4. Within this subset R_i is uniformly distributed between 1 and N_i , where N_i is the number of objects in each subset. With this amendment to the method, the construction of the τ statistic remains the same with τ once again being defined as,

$$\tau = \frac{\sum_i (R_i - E)}{\sqrt{\sum_i V}} \quad (2.62)$$

where the expectation value E and variance V are the as in Equation 2.58. The statistic τ thus gives us a measure of correlation of the data-set whilst properly accounting for the truncation of the survey due to the flux limit.

In Efron and Petrosian (1999) they extend this method for data which has a double truncation i.e. a survey that has distinct faint *and* bright flux-limits. The natural progression was to explore cases where the value of $|\tau| \geq 1$, which implies random variables L and z cannot be considered separable. In this case they assume that there is some form of luminosity evolution in the underlying data. And so by adopting a parametric form for luminosity evolution that varies with an evolutionary parameter, k , they were able to show that for a particular value of k , $\tau(k)$ would equal 0.

Maloney and Petrosian (1999) apply these methods on various quasar samples to determine the density functions and the luminosity evolution. These techniques have also been used by Hao and Strauss (2005) for AGN's to test the correlation between the host galaxy and nuclear luminosity.

2.5 Conclusions

Throughout this chapter we have examined all the innovative statistics developed for the study of the LF. Figure 2.5 shows a time-line diagram which charts the genealogical progress of all these statistics. We have shown that whilst the traditional number count *classical* approach is straightforward in its construction, it is limited by its assumption of spatial homogeneity - a limitation also shared with the $1/V_{\max}$ estimator. Nevertheless, this seems to have not deterred observers from applying Schmidt's estimator, as is evident from its numerous extensions to multiple surveys with varying flux limits, diameter-limited surveys, fitting generic LFs and adaptation to photometric redshifts and so on. It should be noted, however, that despite all these extensions to the basic estimator, the predictive power of V/V_{\max} remains limited. It is difficult to determine whether a significant departure from the expectation value of V/V_{\max} ($\approx 1/2$) is indeed due to effects such as clustering and/or evolution or simply an indication of underlying incompleteness of the catalogue. Despite this, $1/V_{\max}$ and V/V_{\max} remain one of the most widely applied non-parametric statistical tools for constructing LFs and testing for evolution respectively.

Although, the construction of the ϕ/Φ estimator offered a way to effectively circumvent the assumption of homogeneity by the cancellation of the density terms, it is still a statistic that required the binning of data. Later extensions to the method have rendered it a parametric one where an analytical form of the LF is assumed.

The C^- method could be considered as a breakthrough since no binning of the data was required with no dependancy on spatial distribution. And yet, despite the fact that it was pioneered just three years after Schmidt's estimator and eight years before ϕ/Φ , it did not grow in popularity as other tests did. Moreover, as Petrosian (1992) demonstrated, all non-parametric methods are essentially variations of the C^- method. It is therefore a slight mystery as to why this approach has not been applied more often. Perhaps this is due to other maximum likelihood estimators (MLE) that have developed as a result.

Probably the most notable of the MLE's is the so called STY estimator named after its developers, Sandage, Tammann and Yahil in 1979. Its main appeal is that one must adopt a analytical form for the LF and estimate its parameters via a maximum likelihood process. A common form to adopt is the Schechter function (Schechter, 1976). Approaching the problem this way avoids many of the problems associated with binning techniques and also avoids issues of density inhomogeneities. Efstathiou et al. (1988) extended STY's idea to the non-parametric case by replacing the analytical

form of the LF with a series of step functions. This stepwise maximum likelihood method (SWML), often referred to as the EEP method, has become the most favoured estimator in recent times.

We have also briefly considered more recent sophisticated estimators that have emerged within the last two years. The first one we looked at, developed by [Schafer \(2007\)](#), states key points for its construction that include making no assumption of independence between redshift and absolute magnitude and incorporating a varying selection function into the method. The second by [Kelly et al. \(2008\)](#), takes a Bayesian approach to constructing the LF, where the LF is approximated as a mixture of Gaussian functions. Whether approaches such as these become popular remains to be seen, since shifts towards seemingly more convoluted techniques can take time to catch on. The fact that $1/V_{\text{max}}$ is still very much in use today is an indication that perhaps simpler is better?

Finally, we took a step back and turned our attention to one of the basic assumptions of all non-parametric estimators - the assumption of separability between the luminosity function ϕ and the density function $\rho(z)$. By building on the C^- method, [Efron and Petrosian \(1992\)](#) constructed test statistic that can serve as a test of correlation between the assumed independent variables M and z and furthermore be used to constrain pure luminosity models where any such evolution would introduce a correlation in the (M, z) distribution. This work was the foundation for a new completeness test introduced by [Rauzy \(2001\)](#) which we describe in detail in the following chapter and is the basis for the work presented thereafter in this thesis.

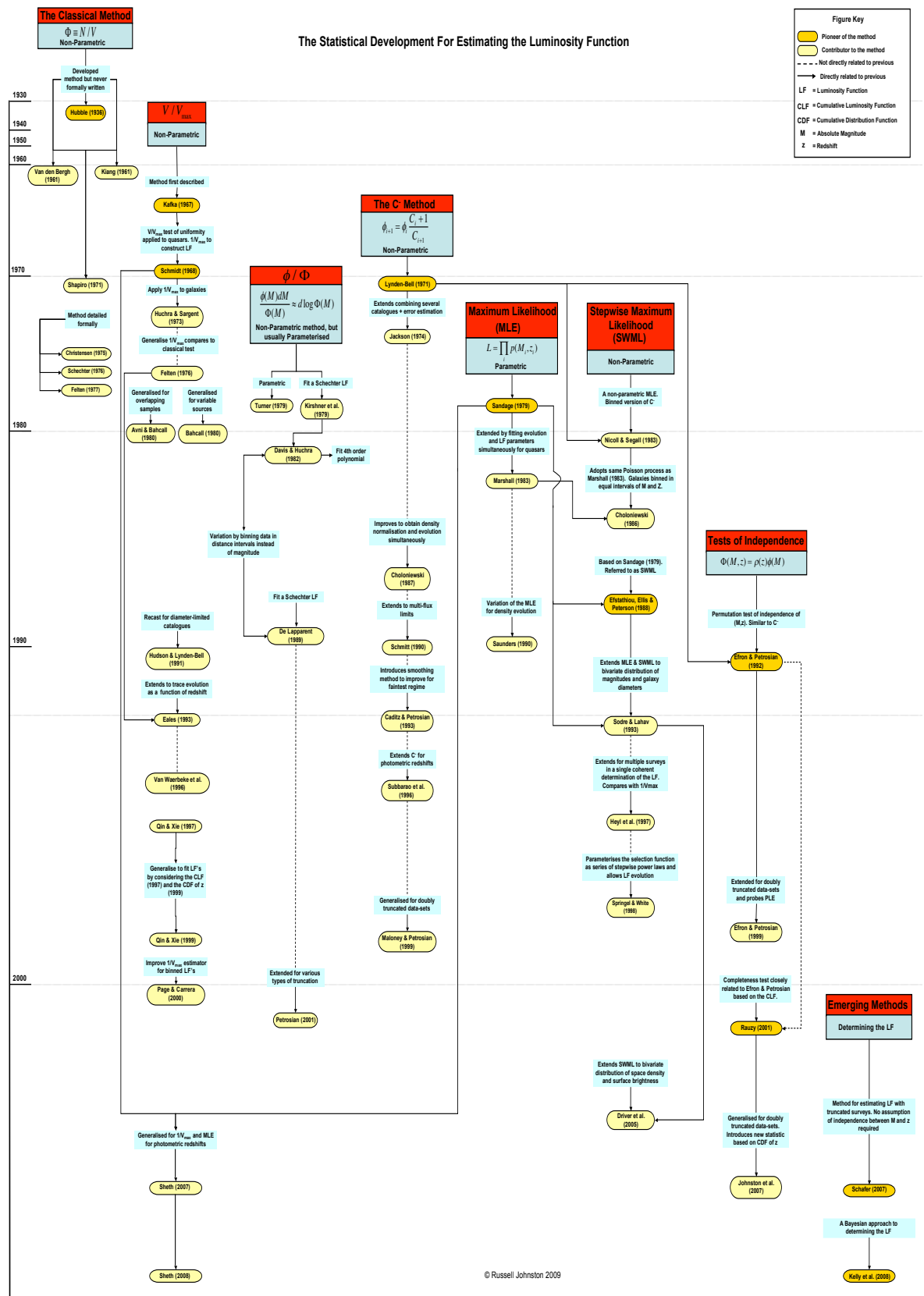


Figure 2.5: Schematic charting the development of all the major statistical methods that estimate the galaxy LF.

Chapter 3

Review of the Rauzy *ROBUST* Method of Completeness

“Do not put your faith in what statistics say until you have carefully considered what they do not say.”

William W. Watt

This chapter reviews the original *ROBUST* test for completeness which was introduced by Stéphane Rauzy (2001; hereafter R01) and which represents the foundation of the research presented in this thesis. We will also demonstrate through analysis of the recent Millennium Galaxy Catalogue (MGC) redshift survey how this test is applied and how the results are interpreted.

3.1 Constructing the R01 Completeness Statistic

The motivation for the original development of this statistical method was to introduce a non-parametric test that could determine the *true* completeness limit in apparent magnitude of a magnitude-redshift survey whilst retaining as few model assumptions as possible.

3.1.1 Assumptions and statistical model

The fundamental assumption of the Rauzy completeness test is that the luminosity function of a galaxy population is universal (does not depend on the 3-D redshift space position $\mathbf{z} = (z, l, b)$ of the galaxies or their type). Although this assumption is

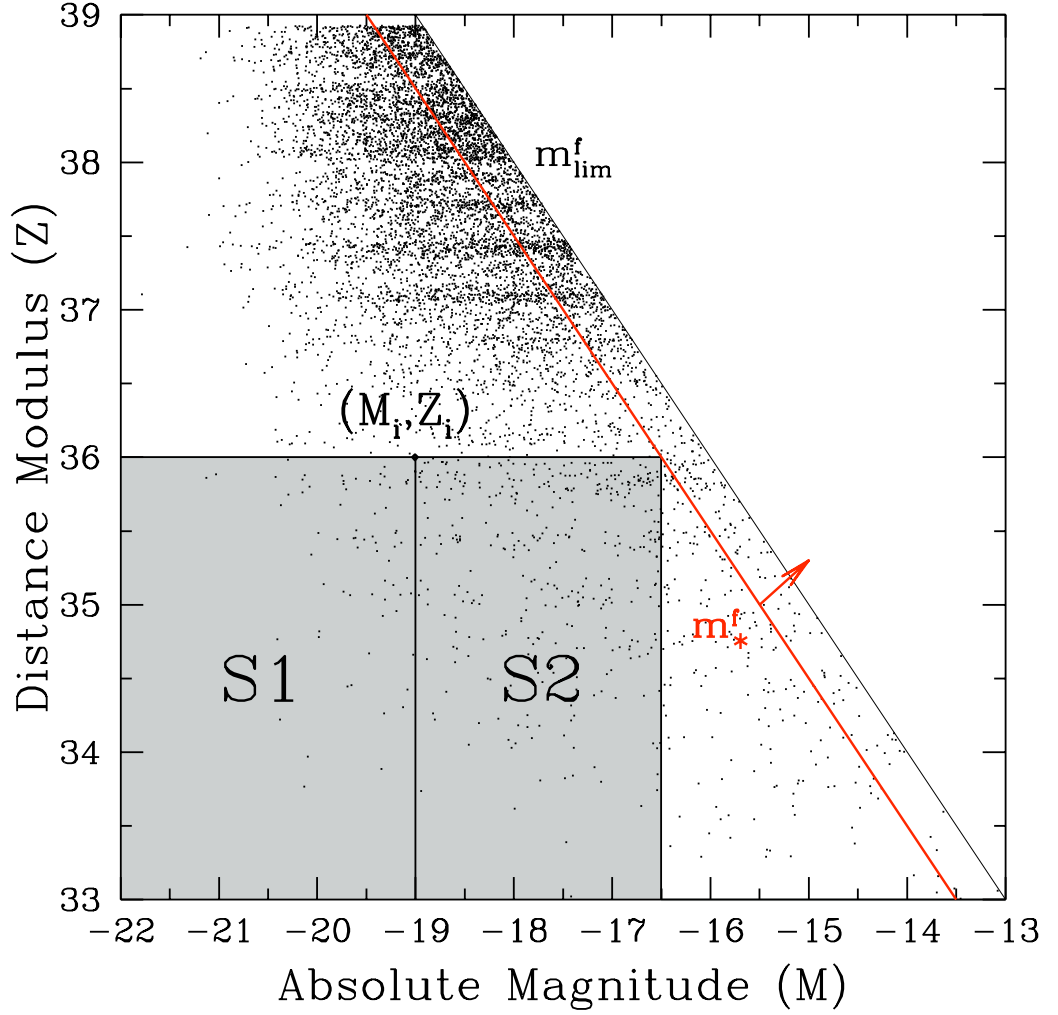


Figure 3.1: Schematic diagram illustrating the construction of the rectangular regions S_1 and S_2 , defined for a typical galaxy at (M_i, Z_i) . The plot shows the original R01 construction of regions S_1 and S_2 which are defined for ‘trial’ faint apparent magnitude limit m_*^f . Also shown is the true faint apparent magnitude limit m_{lim}^f , within which the rectangular regions S_1 and S_2 contain a joint distribution of M and Z that is separable.

restrictive, note that it is common to most classical number counts tests of completeness and indeed (when applied in the context of assessing magnitude completeness) to the [Schmidt \(1968\)](#) V/V_{max} test too as previously discussed in Chapter 2. Moreover, the results derived in Appendix 1 of [Rauzy et al. \(2001\)](#) imply that the completeness test of R01 remains valid in the case of pure density evolution. Note also that the completeness test remains valid for the case of pure luminosity evolution provided that the correct evolutionary model is applied to account for the (assumed known) functional dependence of the mean luminosity at given redshift. Following R01, we introduce the

corrected distance modulus Z , defined as

$$Z_{\text{corr}} = 5 \log_{10} \left(\frac{d_L}{1 \text{ Mpc}} \right) + 25 + k(z) + e(z) + A_g(l, b), \quad (3.1)$$

where, $k(z)$ and $e(z)$ are k -corrections and evolutionary corrections respectively and $A_g(l, b)$ is an extinction correction dependent on galactic co-ordinates (l, b) . The luminosity distance d_L , in a Friedmann Universe, is defined as

$$d_{L_i} = (1 + z_i) \left(\frac{c}{H_0} \right) \int_0^{z_i} \frac{dz}{\sqrt{(1+z)^3 \Omega_{m0} + \Omega_{\Lambda0}}}, \quad (3.2)$$

where Ω_{m0} and $\Omega_{\Lambda0}$ represent the present-day dimensionless matter density and cosmological energy density constant respectively, z_i is the redshift of the i^{th} galaxy in the survey, c is the speed of light, and H_0 is the Hubble constant.

R01 goes on to define the corrected absolute magnitude, M_{corr} as:

$$M_{\text{corr}} = m - 5 \log_{10}(d_L) - 25 - k(z) - e(z) - A_g(l, b). \quad (3.3)$$

Neglecting for the moment any observational selection effects, the joint probability density in position and absolute magnitude for the galaxy population can be written as

$$dP_{zM} \propto dP_z \times dP_M = \rho(z, l, b) dl db dz \times f(M) dM, \quad (3.4)$$

where $\rho(z, l, b)$ is the 3D redshift space distribution function of the sources along the past light-cone and $f(M)$ is the galaxy luminosity function, defined following e.g. [Binggeli et al. \(1988\)](#). We now take as our null hypothesis that the selection effects are separable in position and apparent magnitude, and that the observed sample is complete in apparent magnitude for those objects which are brighter than a specified faint apparent magnitude limit, m_{lim}^f . Under this null hypothesis the selection function $\psi(m, z, l, b)$ can be written as

$$\psi(m, z, l, b) \equiv \theta(m_{\text{lim}}^f - m) \times \phi(z, l, b), \quad (3.5)$$

where $\theta(x)$ is the Heaviside or ‘step’ function defined as

$$\theta(x) = \begin{cases} 1 & \text{if } x \geq 0, \\ 0 & \text{if } x < 0, \end{cases} \quad (3.6)$$

and $\phi(z, l, b)$ describes the selection effects in angular position and observed redshift. Taking into account this model for the selection effects, the probability density function describing the joint distribution of absolute magnitude M and corrected distance

modulus Z for the observable galaxy population may therefore be written as

$$dP = \bar{h}(Z)dZ f(M)dM \theta(m_{\text{lim}}^f - m), \quad (3.7)$$

where $\bar{h}(Z)$ is the probability density function of Z for observable galaxies, marginalised over direction on the sky, *i.e.*

$$\bar{h}(Z) = \int_l \int_b h(Z, l, b) dl db, \quad (3.8)$$

and the integrand $h(Z, l, b)$ is equal to the (suitably normalised) product of the 3-D redshift space density $\rho(z, l, b)$ and the selection function $\phi(z, l, b)$, re-expressed as a function of Z rather than redshift, z .

3.1.2 Defining and estimating the random variable ζ

Note from Equation 3.7 that the faint apparent magnitude limit introduces a correlation between the variables M and Z for observable galaxies. To retain the assumption of separability, the key to the construction of the Rauzy completeness statistic is the definition of the random variable ζ . Since, by the introduction of a faint limit we have a correlation between M and Z , we construct the random variable, ζ , in such a way as to render M and Z separable again. This is achieved by constructing a rectangular region within m_{lim}^f (illustrated in Figure 3.1) and defining within that region the variable ζ satisfying:

$$\zeta = \frac{F(M)}{F[M_{\text{lim}}(Z)]}, \quad (3.9)$$

where $F(M)$ is the Cumulative Luminosity Function (CLF) defined as,

$$F(M) = \int_{-\infty}^M f(x) dx. \quad (3.10)$$

It follows immediately from its definition that ζ has the property of being uniformly distributed between 0 and 1. Under the null hypothesis as described in section 3.1.1, the random variable ζ can be estimated directly from the data, without prior knowledge of the functional form of the CLF.

Consider a galaxy located at (M_i, Z_i) as illustrated in Figure 3.1. The galaxy defines the rectangular region $S_1 \cup S_2$ within a trial magnitude limit, m_{*}^f . The regions, S_1 and S_2 are defined as:

- $S_1 = \{M \leq M_i \text{ and } Z \leq Z_i\}$

- $S_2 = \{M_i < M \leq M_{\text{lim}} \text{ and } Z \leq Z_i\}$

We then define the number of galaxies, r_i , contained within S_1 as:

$$\frac{r_i}{N_{\text{gal}}} = F(M_i) \times \frac{1}{A} \int_{-\infty}^{Z_i} h(Z, l, b) dZ dl db \quad (3.11)$$

where N_{gal} is the number of galaxies in the dataset. Similarly, the number of galaxies, n_i , contained within the region $S_1 \cup S_2$ is defined as:

$$\frac{n_i}{N_{\text{gal}}} = F[M_{\text{lim}}(Z_i)] \times \frac{1}{A} \int_{-\infty}^{Z_i} h(Z, l, b) dZ dl db \quad (3.12)$$

Therefore, for an unbiased estimate (see e.g. [Efron and Petrosian, 1992](#)) of ζ we have,

$$\hat{\zeta}_i = \frac{r_i}{n_i + 1}. \quad (3.13)$$

which has an expectation value E_i and variance V_i given by, respectively,

$$E_i = E(\hat{\zeta}_i) = \frac{1}{2}, \quad V_i = E\left[\left(\hat{\zeta}_i - E_i\right)^2\right] = \frac{1}{12} \frac{n_i - 1}{n_i + 1}. \quad (3.14)$$

Note that V_i tends towards the variance of a continuous uniform distribution between 0 and 1 when n_i is large.

We can, therefore, combine the estimator $\hat{\zeta}_i$ for each observed galaxy into a single statistic, T_c , which we can use to test the magnitude completeness of our sample for adopted trial magnitude limits m_*^f . T_c is defined as,

$$T_c = \sum_{i=1}^{N_{\text{gal}}} \left(\hat{\zeta}_i - \frac{1}{2}\right) / \left(\sum_{i=1}^{N_{\text{gal}}} V_i\right)^{\frac{1}{2}}. \quad (3.15)$$

If the sample is complete in apparent magnitude, for a given trial magnitude limit, then T_c should be normally distributed with mean zero and variance unity. If, on the other hand, the trial faint magnitude limit is fainter than the true limit, T_c will become systematically negative, due to the systematic departure of the $\hat{\zeta}_i$ distribution from uniformity on the interval $[0, 1]$.

3.1.3 Implementing the T_c statistic

Following R01, we calculate T_c as a function of trial faint magnitude limit denoted as m_* . Slicing the data in this way allows us to determine a differential measure of completeness over the whole range of apparent magnitudes.

Figure 3.2 illustrates the systematic way in which we can computationally exclude regions in the M - Z plane to estimate $\hat{\zeta}_i$ for each galaxy located at (M_i, Z_i) . For an initial value of m_* we remove all galaxies $m_i > m_*$. From this we calculate an absolute magnitude limit for each galaxy such that,

$$M_{\text{lim}}(Z_i) = m_*^i - Z_i, \quad (3.16)$$

and remove all galaxies fainter than $M_{\text{lim}}(Z_i)$. Finally, before we estimate $\hat{\zeta}_i$ we can remove the remaining galaxies that lie at a distance modulus greater than Z_i . By applying these initial steps we improve the efficiency of the T_c calculation. Thus, for each $T_c(m_*)$ we now estimate the random variable $\hat{\zeta}_i$ for each galaxy at (M_i, Z_i) by simply counting the number of galaxies in the pre-determined regions, r_i and n_i (given by equations 3.11 and 3.12 respectively) whilst simultaneously calculating the variance, V_i .

We then combine all the $\hat{\zeta}_i$'s for our m_* and calculate a value for the completeness statistic, $T_c(m_*)$ from equation 3.15. Therefore, for a sample complete in apparent magnitude, one expects the value of $T_c(m_*)$ to be normally distributed around zero with sampling fluctuations of dispersion of unity. However, as Figure 3.3 illustrates, when the trial magnitude limit, m_* moves beyond the *true* magnitude limit of the survey, m_{lim} , the number of galaxies contained n_i contained in $S_1 \cup S_2$ drops considerably (see the blue region in Figure 3.3) the result of which is a sharp, systematic drop in the value of T_c statistic indicating the limit of the survey.

As is detailed in R01, these confidence levels of rejection for $|T_c| < 1$, $|T_c| < 2$, $|T_c| < 3$ are 84.13%, 97.72% and 99.38% respectively, and these confidence levels may be then used to make a choice of the appropriate completeness limit. These numbers correspond to probabilities for a normal distribution with a one sided rejection test. We have chosen a significance level of 3σ throughout this thesis as our adopted criterion for identifying completeness limits.

3.2 Applying T_c to the Millennium Galaxy Catalogue (MGC)

In this section we now apply the T_c test statistic to the Millennium Galaxy Catalogue (MGC) to carefully illustrate how the results are interpreted.

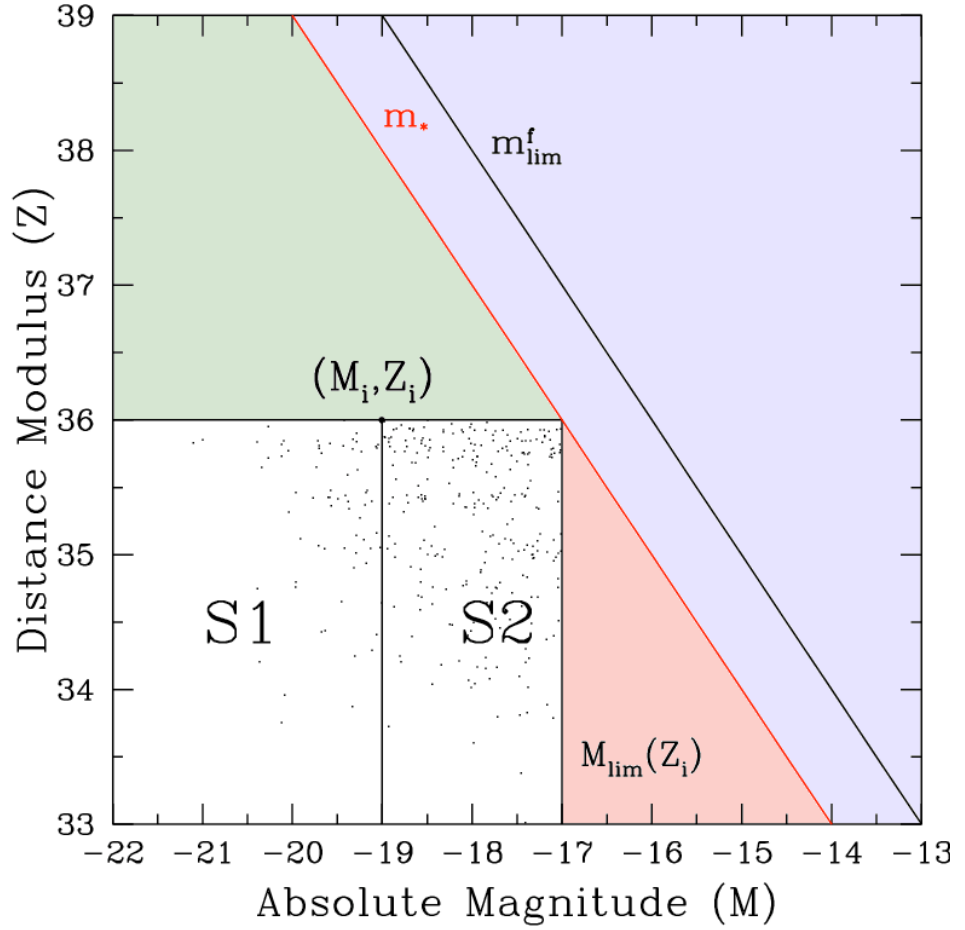


Figure 3.2: Schematic diagram illustrating the way in which we computationally select the S_1 and S_2 regions, defined for a typical galaxy at (M_i, Z_i) . The plot shows the original R01 construction of regions S_1 and S_2 defined from the i^{th} galaxy which are defined for ‘trial’ apparent magnitude limit m_* shown as the red diagonal line. Also shown is the true faint apparent magnitude limit m_{lim}^f (black diagonal line), within which the rectangular regions S_1 and S_2 contain a joint distribution of M and Z that is separable. The purple region shows how we firstly remove all galaxies m fainter than m_i^* . The red triangular shaded region represents the galaxies that are then removed which are fainter than $M_{\text{lim}}(Z_i)$. Finally, we remove the remaining galaxies that lie at a distance moduli greater than Z_i . Processing the data-set in this way allowed us to compute the remaining S_1 and S_2 regions with far greater efficiency.

3.2.1 The Data

The Millennium Galaxy Catalogue (MGC) is a medium-deep, B -band imaging survey, spanning 30.9 deg^2 that is fully contained within the 2dFGRS (Colless, 2001) and SDSS-DR1 (Abazajian, 2003). As of 2005 the full catalogue contained 10095 galaxies out to published limiting apparent magnitudes of $13.0 < m_{\text{lim}} < 20.0 \text{ mag}$ - see e.g. Cross et al. (2004) for more details. The photometry was obtained with the Wide Field Camera

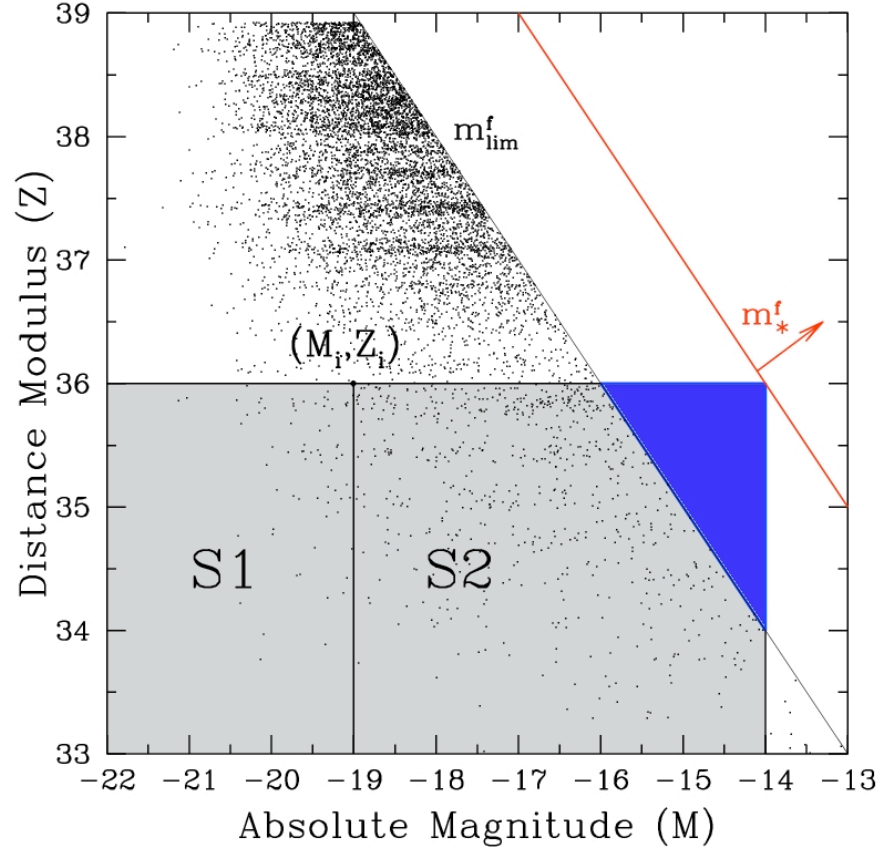


Figure 3.3: Diagram illustrating movement of m_*^f beyond the apparent magnitude limit of a given survey. As this diagram shows, as the trial magnitude limit m_*^f moves beyond the apparent magnitude of the survey, the number of galaxies in region S_2 drops significantly as indicated by the blue region. Therefore, the resulting T_c curve drops sharply below 3σ as m_*^f moves beyond a sharp limit, indicating the *true* limit of the dataset.

on the 2.5 m Isaac Newton Telescope in La Palma. The spectroscopy was constructed mainly from the redshifts obtained in the Two Degree Field Galaxy Redshift Survey (2dFGRS) and the Sloan Digital Sky Survey (SDSS). Additional redshifts were obtained from 2QZ (Croom et al., 2001), the Nasa Extragalactic Database (NED), Francis et al. (2004) and Impey et al. (1996) (LSBG). In addition, the MGC team measured their own redshifts using the spectrograph on the Anglo-Australian Telescope for galaxies in the catalogue that had no assigned redshift.

3.2.2 Selection Limits & Cosmology

The MGC team adopted the same redshift quality procedure that was originally applied to the 2dFGRS. Galaxies were assigned a redshift quality, Q_z , from 1 to 5: $Q_z = 1$ indicates that no redshift could be estimated; $Q_z = 2$ indicates a possible redshift measurement; $Q_z = 3$ is a ‘probable’ redshift with 90% confidence; $Q_z = 4$ is a reliable redshift with 99% confidence; and $Q_z = 5$ is a reliable redshift with a high-quality spectrum. Therefore, galaxies with a published redshift quality $Q_z \geq 3$ were selected. For ease of comparison we have imposed the same redshift selection as [Driver et al. \(2005\)](#) (D05) – i.e only galaxies in the range $0.013 < z < 0.18$ were included. From the parent catalogue of 10,095 galaxies this yields a sub set of 7,878 galaxies.

In keeping with D05 and for ease of comparison we adopt the same cosmology with present-day dimensionless matter density parameter, $\Omega_{m0} = 0.3$ and cosmological constant term $\Omega_{\Lambda0} = 0.7$. For this survey we have adopted a value of $H_0 = 100 \text{ kms}^{-1} \text{ Mpc}^{-1}$ for the Hubble constant.

3.2.3 k - and evolutionary corrections

Where appropriate we have applied the k -corrections and evolutionary corrections as described in greater detail in D05. Such corrections applied to the MGC catalogue are galaxy specific and are determined by using a galaxy spectrum template-fitting technique for each individual galaxy. Corrections for evolution within the data-set were characterised by a global correction of the form,

$$L = L_0(1 + z_i)^\beta \quad (3.17)$$

where L is the luminosity, z is the redshift and β is an evolution parameter. To convert to absolute magnitudes we know the relation,

$$\left(\frac{L}{L_0}\right) = 10^{-0.4(M-M_0)} \quad (3.18)$$

giving the evolution correction in the form,

$$E(z) = -\beta \times 2.5 \log_{10}(1 + z_i). \quad (3.19)$$

The value of β was determined in [Driver et al. \(2005\)](#) to be in the range $-2.0 < \beta < 1.25$ and a global value of $\beta = 0.75$ was considered suitable.

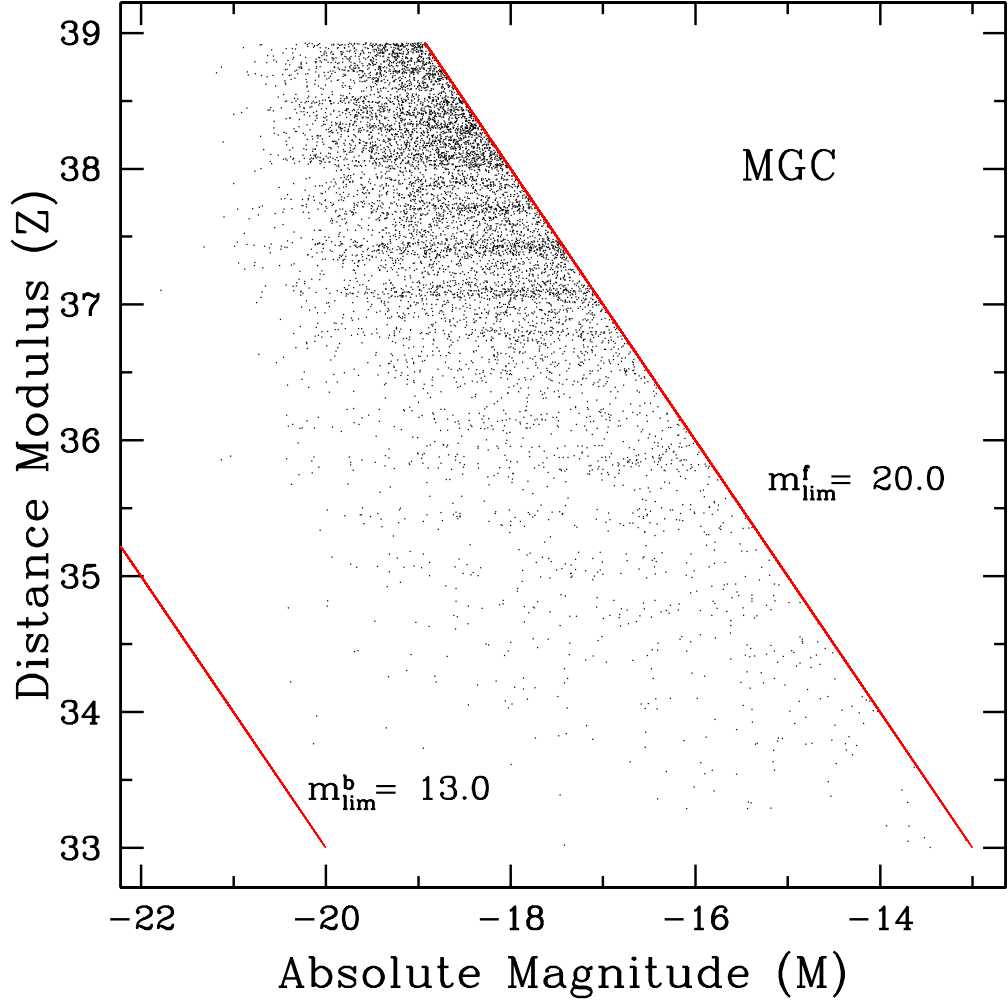


Figure 3.4: A subset of the MGC M - Z distribution. In this sample we have 7878 galaxies from parent catalogue of 10,095 galaxies. The red diagonal lines indicate the survey apparent magnitude limits. Since we are considering a subset we note that there are no galaxies present at the published bright apparent magnitude limit $m_{\text{lim}}^b = 13.0$ mag

3.2.4 Results

Figure 3.4 shows the M - Z distribution of the selected data-set. It should be noted that although MGC has a published bright limiting magnitude of $m_{\text{lim}}^b = 13.0$ mag, one can clearly see from Figure 3.4 that there is no clear sharp bright cut-off, an effect that will play a crucial role in subsequent chapters.

By applying T_c we confirm the completeness of the MGC data demonstrated in Figure 3.5 where the dashed curve shows the T_c statistic, as a function of trial magnitude limit, computed using apparent magnitudes that have *not* been $(k+e)$ -corrected,

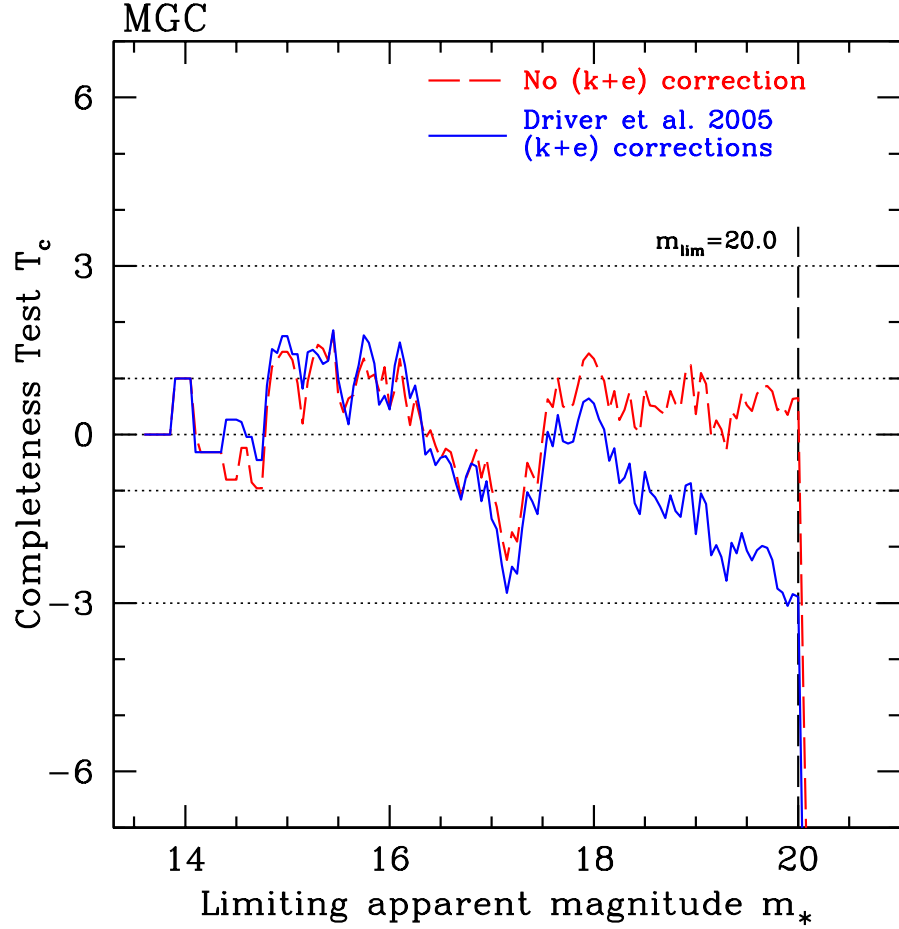


Figure 3.5: Results for the completeness statistic T_c applied to MGC without explicitly accounting for a bright apparent magnitude limit. The red dashed line represents the results from the MGC data-set without any k or evolution corrections applied to the dataset. The blue solid line is the result we have apply the k -corrections and evolution corrections with $\beta = 0.75$. The decay in the curve for the corrected magnitudes is thought to be due to the mixing of different galaxy populations as a result of individual corrections, resulting in a ‘fuzzy’ magnitude limit close to $m_{\text{lim}} = 20.0$ mag.

but have been corrected for extinction only. We have chosen to compute T_c for an m_* moving to increasingly fainter magnitudes in increments of 0.1 from the brightest galaxy in the set, $m = 13.63$ mag. The figure clearly shows that the T_c statistic remains within the 3σ limits – consistent with being complete in apparent magnitude – up to the published magnitude limit of 20.0 mag, and then drops very sharply for trial apparent magnitude limits beyond 20.0 mag. To further illustrate this point we can observe the (ζ, Z) distribution for two different slices of m_* shown in Figure 3.6. By definition, the ζ for a complete sample at any given m_* will be uniformly distributed between 0 and 1. The left hand panel in the figure considers this distribution at the

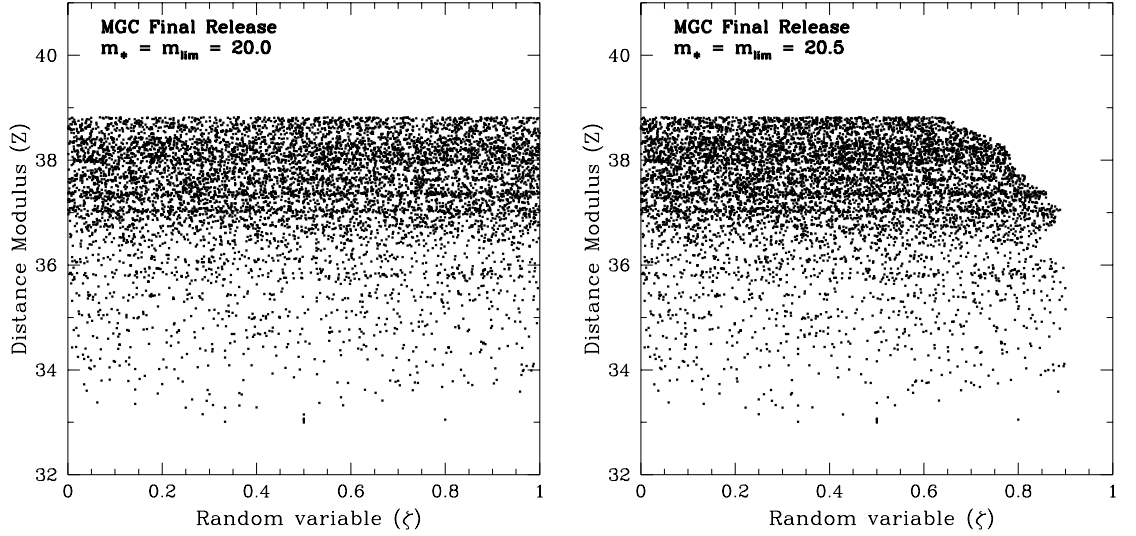


Figure 3.6: The random variable ζ vs distance modulus, Z , for MGC at different values of m_* . The left hand panel shows the (ζ, Z) distribution for $m_* = 20.0$ mag - the published magnitude limit of the survey. The T_c curve in Figure 3.5 showed the data-set was indeed complete up to m_{lim} which is also verified by this plot since the left hand panel shows a uniform distribution of ζ between $[0,1]$. The sharp drop in T_c that was observed as we moved beyond the apparent magnitude limit is seen just as dramatically in the (ζ, Z) distribution as shown in the right hand panel which is for $m_* = 20.5$ mag. Here we observe that the distribution of points are no longer uniform between $[0,1]$ particularly at the top right section, due to the deficit in the number of galaxies in the $S2$ region, as already shown in Figure 3.3. This in turn renders the $S1$ and $S2$ regions un-separable which manifests as a systematic shift of the points in the (ζ, Z) plane.

magnitude limit of the survey, $m_{\text{lim}} = 20.0$ mag. As we have already stated, the T_c results indicate the survey is indeed complete up to and including this limit, and the distribution on the left and panel also confirms this. It is clear by eye that the ζ is uniform on the interval $[0,1]$. On the right hand panel of Figure 3.6 we now consider the (ζ, Z) distribution for an $m_* = 20.5$ which takes us beyond the sharp cut-off of the survey and equates to a value of $T_c \approx -32.3\sigma$. What results is a curving of the (ζ, Z) distribution at the top right corner. If we now consider the solid blue curve on Figure 3.5, we once again see the T_c statistic as a function of trial magnitude limit but now computed for $(k+e)$ -corrected apparent magnitudes. Although the MGC dataset is still consistent with being complete up to the published magnitude limit of 20.0 mag, there is a noticeable departure in the behaviour of T_c from that for the uncorrected dataset: for trial magnitude limits in the range $18.0 \lesssim m_* \leq 20.00$ mag, T_c for the corrected dataset exhibits a slow decline, before again dropping sharply beyond 20.0

mag. The most likely explanation for this feature seems to lie in the way in which the dataset is selected and corrected. The raw dataset, with uncorrected magnitudes, has the same magnitude limit imposed on all galaxies independent of their galaxy type. If, then, each galaxy is individually k -corrected, the resultant overall magnitude limit for the corrected data will become ‘fuzzy’ without a sharp cut-off. Furthermore, different galaxy populations will be scattered differently, leading to a smooth decrease close to the original uncorrected magnitude limit. This effect will be compounded slightly by the global e -correction which is also incorporated. On the other hand, if we do *not* apply the $(k+e)$ -corrections, the original magnitude limit remains well defined (albeit now without explicitly accounting for the effects of evolution and redshifting of each galaxy’s spectral energy distribution). Therefore, to obtain an improved measure of completeness which *does* properly incorporate $(k+e)$ -corrections, one could apply R01 separately to subsets of different galaxy type. This would, in principle, lead to the definition of different apparent magnitude limits for different galaxy types. One could of course apply type-dependent corrections and impose an apparent magnitude cut-off on the whole data-set based on the T_c results derived from m_{corr} . In any event, it is clear from Figure 3.5 that the impact on the inferred ‘global’ apparent magnitude limit of applying, or not, $(k+e)$ -corrections to the MGC dataset is small.

3.3 Conclusions

In this chapter we have reviewed the construction of the original Rauzy test for completeness, T_c and shown how it is applied to real data. For this demonstration we used the MGC redshift survey. It was shown that the magnitude completeness of the survey was confirmed up to its published magnitude limit of $m_{\text{lim}} = 20.0$ mag. Interestingly, however, we found that when we incorporated $(k+e)$ -corrections, a noticeable (although not statistically significant) departure from the results obtained where no corrections were added close to the magnitude limit. A possible cause for this effect could be the mixing of galaxy populations to which k - and global e -correction is then applied – resulting in a slightly *blurred* magnitude limit.

Chapter 4

Extending the Method for Doubly Truncated Surveys

“An unsophisticated forecaster uses statistics as a drunken man uses lamp-posts - for support rather than for illumination.”

Andrew Lang (1844-1912)

This chapter will show in detail the development and extension of the original R01 completeness test (see Chapter 3 and Rauzy, 2001) which arose directly from the completeness results obtained when analysing the Two Degree Field Galaxy Redshift Survey (2dFGRS) and the Sloan Digital Sky Survey - Early Types (SDSS-Early Types).

4.1 Analysis of the 2dFGRS

4.1.1 The Data

The Two Degree Field (2dFGRS) Galaxy Redshift Survey (Colless, 2001) ran from 1997 to 2003 and, at the time of the completion, was the second largest redshift survey next to the Sloan Digital Sky Survey (SDSS), which has now only just reached completion of the second phase. The survey utilised the AAOmega multifibre spectrograph on the Anglo-Australian Telescope capable of measuring redshifts up to 400 objects simultaneously. The corresponding photometry was taken from the APM galaxy catalogue (Maddox et al., 1990) for galaxies brighter than an apparent magnitude of $m_{b_j} = 19.45$ mag. The survey region covered two strips: one $75^\circ \times 10^\circ$ around the north galactic pole and the other $80^\circ \times 15^\circ$ around the southern galactic pole (see Figure 4.1).

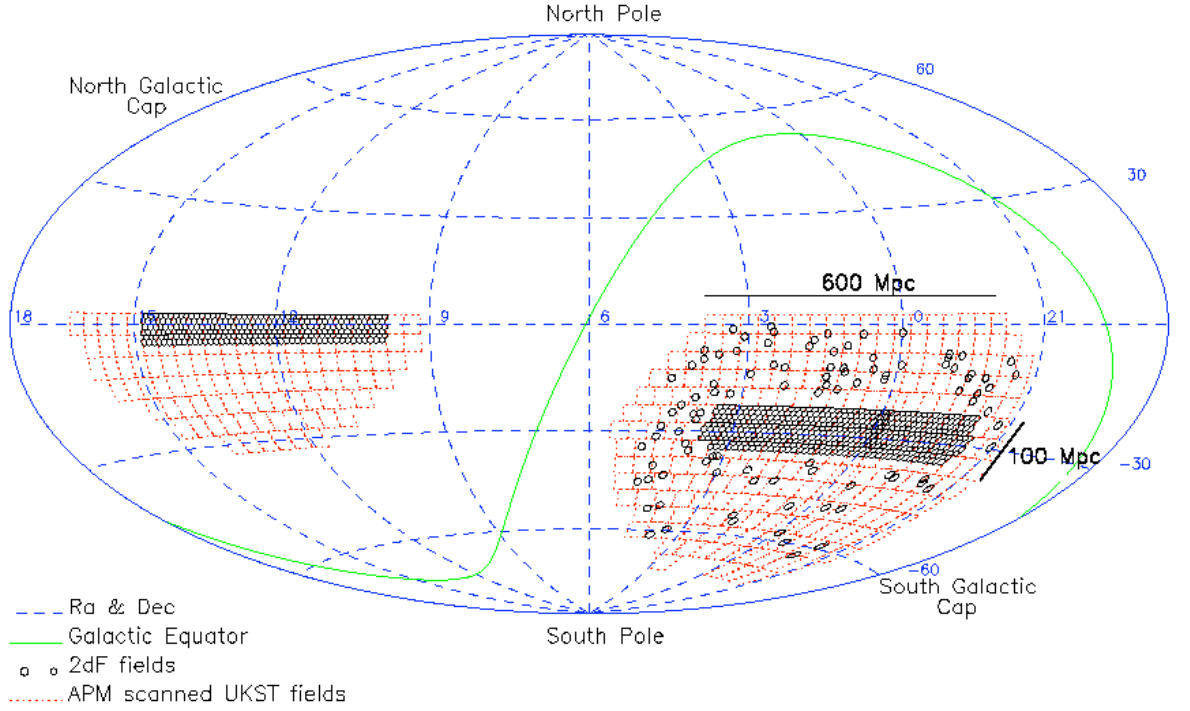


Figure 4.1: Survey map of the 2dFGRS showing the digitised APM plate regions as dashed squares. The selected 2dFGRS regions are shown by the circles. Image courtesy <http://www2.aao.gov.au/2dFGRS/>

There were many objectives for this survey but some of the more crucial ones were to measure: the cosmological mass density (Peacock, 2001), the galaxy power spectrum $P(k)$ on scales up to few hundred Mpc, (Percival, 2001; Tegmark et al., 2002; Outram et al., 2003); the redshift-space distortions of large-scale clustering resulting from peculiar velocity fields; higher order clustering statistics of the galaxy distribution; and of course to provide an extensive spectroscopic database to be used in conjunction with other surveys.

We have used the 2dFGRS public final release data-set, from the ‘*best observations*’ spectroscopic catalogue, which records redshifts for a total of 245,591 sources.

4.1.2 Selection limits

To construct a clean catalogue, we first selected those galaxies with reliable redshifts. The 2dFGRS team classified the redshift for each galaxy with an assigned Quality number from 1 to 5 (Colless, 2001) as we have already detailed in § 3.2.2 on page 60. Therefore, galaxies with a redshift quality of $Q_z \geq 3$ were selected.

We then imposed maximum and minimum limits in redshift following Cross (2001),

$0.015 < z < 0.12$, to minimise the effects of peculiar velocities and isophotal corrections respectively. We used final magnitudes with extinction correction and excluded any galaxies fainter than $m_{b_j} = 19.45$ mag. From the parent catalogue of 245,591 sources we used a total of 110583 galaxies. The corresponding M - Z distribution of our selection is shown in Figure 4.2, with the published survey apparent magnitude limit, $m_{\text{lim}}^f = 19.45$ mag shown as the red diagonal line.

4.1.3 The k - and evolutionary corrections

To explore the effect k - and e - corrections on the T_c statistic we applied various forms that have been used on the 2dFGRS data-set since its release.

The first is a simple global correction, following Driver et al. (1994), and applied in Cross (2001):

$$k(z) = 2.5z \quad (4.1)$$

Similarly, Norberg (2001) used a global correction given by:

$$k(z) = (0.03z)/(0.01 + z^4) \quad (4.2)$$

The final two corrections we considered are type dependent ($k+e$)-corrections defined using a method based on principal-component analysis (PCA), where galaxy type is defined by the η parameter, a linear combination of the first two principal components (Madgwick, 2002). From Norberg (2002b),

$$k(z) = \begin{cases} 2.6z + 4.3z^2 & \text{for } \eta < -1.4 \text{ (early types)} \\ 1.5z + 2.1z^2 & \text{for } \eta > -1.4 \text{ (late types)} \\ 1.9z + 2.7z^2 & \text{(full sample)} \end{cases} \quad (4.3)$$

With further division of type in Norberg (2002a) we have:

$$k(z) = \begin{cases} (2z + 2.8z^2)/(z + 3.8z^3) & \text{for } \eta < -1.4 \\ (0.6z + 2.8z^2)/(1 + 3.8z^3) & \text{for } 1.1 < \eta < -1.4 \\ (z + 3.6z^2)/(1 + 16.6z^3) & \text{for } 1.1 < \eta < 3.5 \\ (1.6z + 3.2z^2)/(1 + 14.6z^3) & \text{for } \eta > 3.5 \end{cases} \quad (4.4)$$

4.1.4 Initial Results

The 2dFGRS revealed some anomalous behaviour in the resulting completeness T_c curves which led us, through a process of experimental investigation, to the conclusion

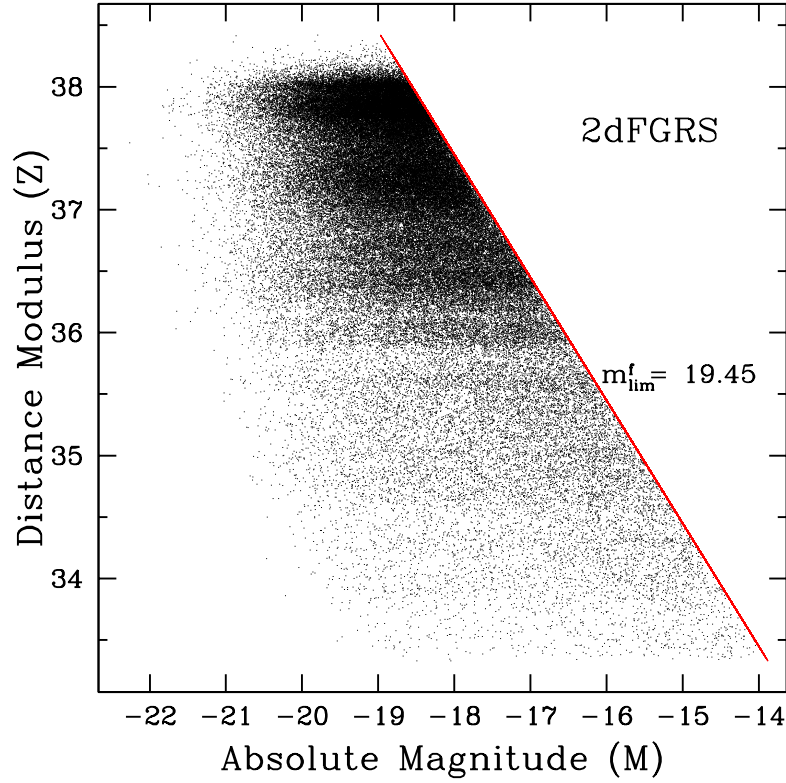


Figure 4.2: The M - Z distribution of our selected sub-set of the 2dFGRS. The red line indicates the published apparent magnitude limit, $m_{\text{lim}}^f = 19.45$ mag.

that this particular survey is inconsistent with magnitude completeness unless one adopts a secondary bright limit. Figure 4.2 shows the distribution of galaxies within the M - Z plane that we are testing. Our initial approach was to apply the original R01 T_c statistic (Chapter 2) to our 2dFGRS selection since the published literature on the survey gives no indication about the presence of a secondary bright limit. Figure 4.3 shows the behaviour of T_c as a function of trial magnitude limit, m_*^f , for five different cases. The solid red curve represents the completeness test with no k - or e -corrections applied. The remaining four curves show T_c with various $(k+e)$ -corrections applied to the 2dFGRS data-set, as explained in the figure key.

If we consider firstly the uncorrected data-set (solid red curve), we see that for $m_* < 14.85$ mag the T_c statistic appears to behave in a manner consistent with what we would expect for a complete sample (although of course this is at the cost of ‘throwing away’ most of the galaxies in the survey by imposing such a low value for the faint limit). However, for higher values of m_*^f the statistic drops dramatically to a minimum value of nearly 8σ below its expectation value of zero at $m_*^f = 16.90$ mag. As we

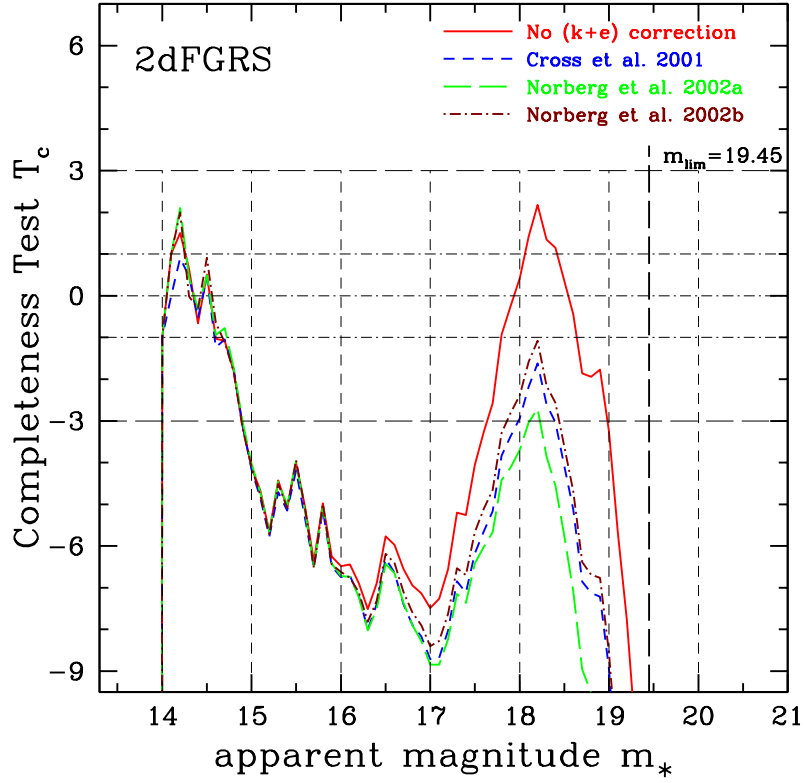


Figure 4.3: The T_c statistic for our entire 2dFGRS sample. The data used for this plot are the combined north and south regions of the survey. The dashed lines represent various $(k+e)$ -corrections that have been applied to this dataset. The striking features indicate a dramatic departure from completeness with a minimum near $m \approx 16.90$ mag and a peak at $m \approx 18.15$ mag. The addition of different types of $(k+e)$ -corrections clearly show that they do not improve the underlying issues that are indicated by T_c .

continue to increase m_*^f , T_c rises sharply to reach a peak at $m_*^f = 18.15$ mag, beyond which the statistic drops dramatically again, exceeding 3σ below its expected value at $m_*^f = 18.60$ mag – i.e. significantly brighter than the published magnitude limit of $m_{\text{lim}}^f = 19.45$ mag.

At first it was thought that the behaviour of T_c could be related to the fact that we have used an uncorrected $(k+e)$ data-set. To address this question consider now the remaining four curves; the dotted and short dashed curves correspond to the Cross (2001) and Norberg (2001) global $(k+e)$ -corrections respectively, whereas the long dashed and solid black curves correspond to the type-dependent $(k+e)$ -corrections of Norberg (2002b,a) as detailed in section 4.1.3. It is clear that the adoption of any of these corrections has very little effect on the completeness statistic compared with the uncorrected case. Indeed, if anything, the addition of $(k+e)$ -corrections appears to yield

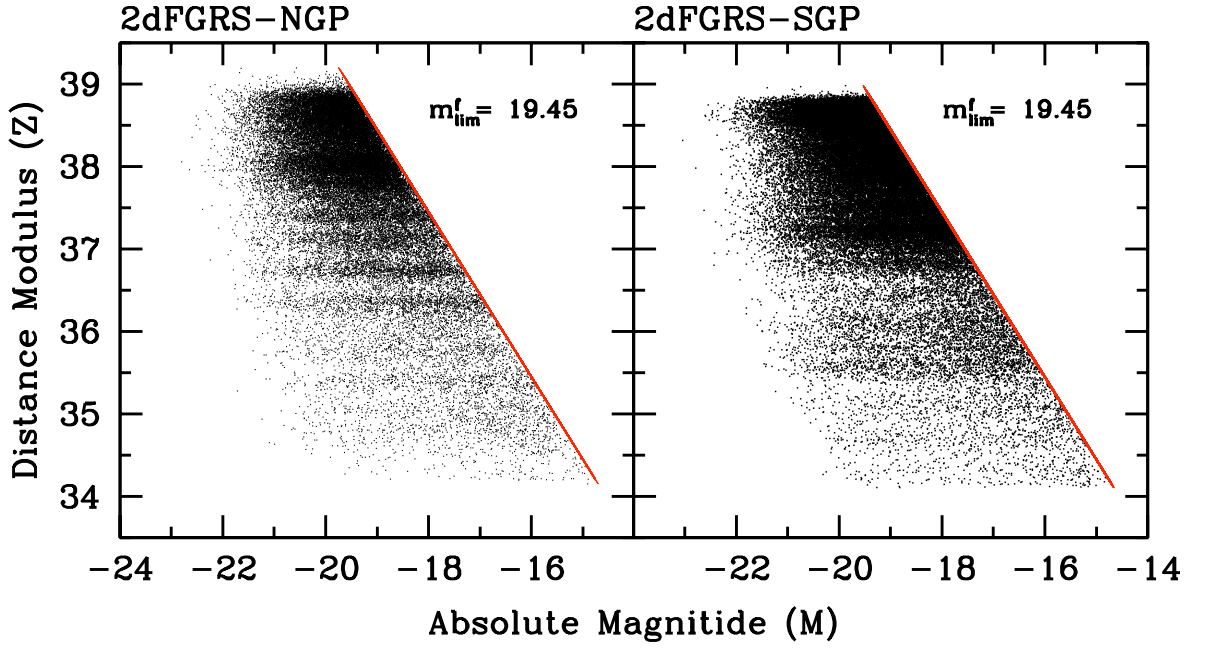


Figure 4.4: The M - Z distribution of the 2dF - Northern and Southern strips.

a T_c statistic which is more strongly inconsistent with a complete sample. This latter effect was discussed in greater detail when we analysed the MGC in section 3.2. It should be noted that the type-dependent ($k+e$)-corrections do not appear to perform significantly better than their global counterparts.

The fact that the value of T_c differs from zero at many standard deviations over a wide range of trial faint magnitude limits is clear evidence that the distribution of M and Z for observable galaxies is *not* separable with these faint limits alone. The physical reason for this is not immediately clear. However, it was thought initially that the cause of incompleteness may lie within one of the regions of the survey. By simply splitting the data-set into the north and south regions (see Figure 4.4) we can show, however, that the T_c results for both regions (Figure 4.5) exhibit the same characteristics as seen in Figure 4.3 for the whole survey.

Our further attempts to determine the cause of incompleteness and, therefore, what had broken the separability between M and Z , led us to investigate the individual photometric plates that make up the survey. The photometric data used in this survey was taken from a subsample of the Automated Plate Measuring-machine (APM) galaxy catalogue. Figure 4.6 shows the magnitude limit mask plate regions for for both, North and South. The colour gradient indicates the varying magnitude limit for each plate which is defined by the change in the photometric calibration for each UKST

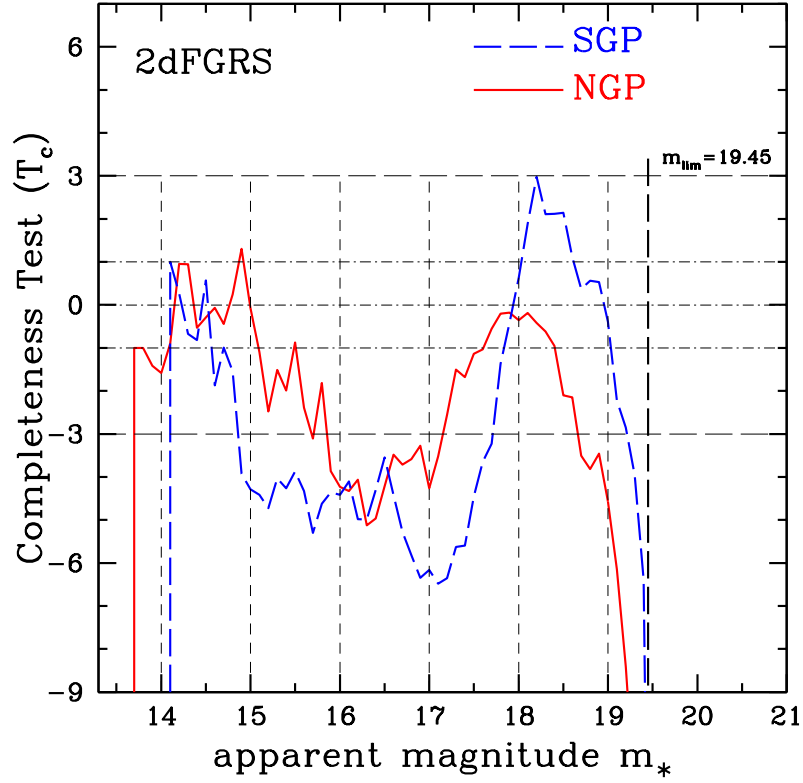


Figure 4.5: The 2dFGRS T_c results for the Northern (NGP) and Southern (SGP) regions. We see similar features to that of the combined NGP and SGP shown in Figure 4.3.

photographic plate and change in the dust extinction correction at each position on the sky. It was thought that perhaps there may have been inconsistencies within one or more plates to which T_c was being sensitive.

There are a total of 30 plates that cover the Northern Galactic Pole (NGP) and 46 plates in the Southern Galactic Pole (SGP) region of the 2dF. Tables 4.1 and 4.2 show in more detail a break down of the number of galaxies we are considering in each plate before and after our selection process. We also show the brightest and faintest magnitudes from plate to plate. Figures 4.7 and 4.8 show plots the T_c statistic obtained by analysing individually each plate in the NGP region and Figures 4.9, 4.10 and 4.11 show the corresponding T_c results for the SGP. For the northern region the results indicate that within each plate there seems to be no evidence of incompleteness, which was at first a puzzling result. However, when examining the southern region we noted that three of the plates, 354, 408 and 477 showed similar behaviour to that of the whole combined survey with each having *marginal* dips in the T_c curve below the -3σ level.

It was at this point that it became apparent that each plate had at most a few

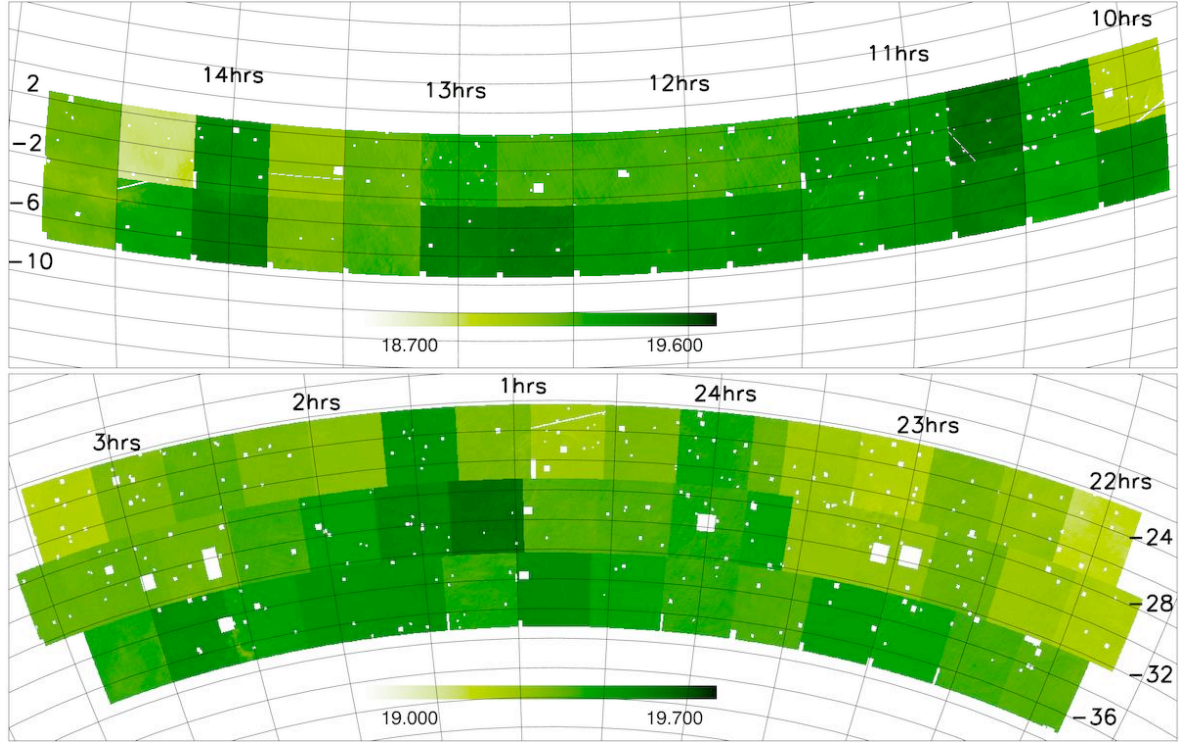
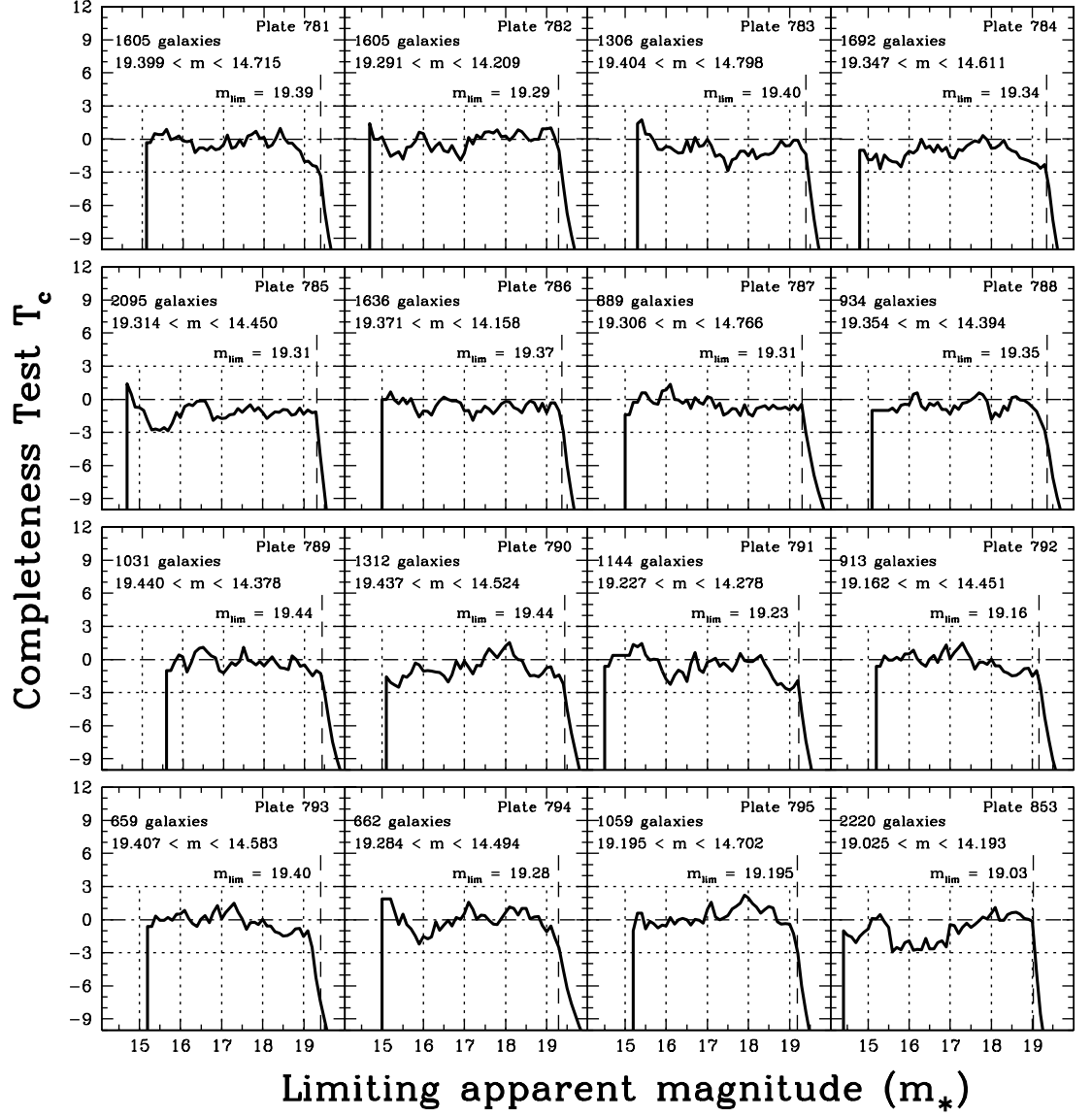


Figure 4.6: The 2dFGRS magnitude limit masks for the NGP (top) and the SGP (bottom). Each square represents one of the APM plates. The variation in colour indicates the varying magnitude limits from plate to plate. Image courtesy of <http://msowww.anu.edu.au/2dFGRS/>

thousand galaxies and by inspection of the corresponding M - Z diagrams it was clear to see that in the majority of plates there was not an obvious, well defined bright limit in apparent magnitude. However, when we looked more closely at those T_c plots that dropped slightly below the -3σ level, namely plates 354, 408 and 477 from Figure 4.12, we could see that each of the corresponding M - Z distributions indicate the presence of a fairly distinct bright apparent magnitude limit. By comparing these to the randomly selected plate 788, where the T_c curve is well behaved (Figure 4.12-top right) we see that on plate 788 there is a slight scatter of nearby bright galaxies that have the effect of blurring any presence of a bright limit. This result led to the conclusion that perhaps the 2dFGRS has indeed a bright limit which should be more carefully accounted for, by incorporating it explicitly into the *ROBUST* method.

Table 4.1: 2dF-Northern Plate Information. This table shows information regarding the Northern Galactic Pole region of 2dF. The *IFIELD* column refers to the UKST plate number. The following two columns show the number of galaxies in each plate before and after our selection criteria respectively. The faintest and brightest galaxies for each plate given in apparent magnitudes are those after selection.

Index	IFIELD	No. in IFIELD (raw)	No. in IFIELD (after selection)	Faintest app. magnitude	Brightest app. magnitude
1	781	3355	1605	19.399	14.715
2	782	2819	1605	19.291	14.209
3	783	2887	1306	19.404	14.798
4	784	3021	1691	19.347	14.611
5	785	3178	2095	19.314	14.450
6	786	3115	1636	19.371	14.158
7	787	1904	889	19.306	14.766
8	788	1806	934	19.354	14.394
9	789	2346	1031	19.440	15.378
10	790	2678	1312	19.437	14.524
11	791	2234	1144	19.227	14.378
12	792	1621	913	19.162	14.451
13	793	1418	659	19.407	14.583
14	794	1359	662	19.284	14.494
15	795	1776	1059	19.195	14.702
16	853	3411	2220	19.025	14.193
17	854	4269	2637	19.323	14.102
18	855	4788	2618	19.450	14.505
19	856	4038	2338	19.335	14.509
20	857	4476	2560	19.360	14.444
21	858	3517	1965	19.239	14.097
22	859	3980	2105	19.225	13.957
23	860	3057	1623	19.249	14.395
24	861	3949	2335	19.241	14.930
25	862	3780	2147	19.266	14.168
26	863	4421	2769	19.217	13.454
27	864	2928	1696	19.106	14.061
28	865	2973	1347	19.386	14.035
29	866	2598	1254	18.933	13.348
30	867	1930	841	19.192	14.570

Figure 4.7: T_c curves for the 2dF-NGP plate numbers: 781 to 853.

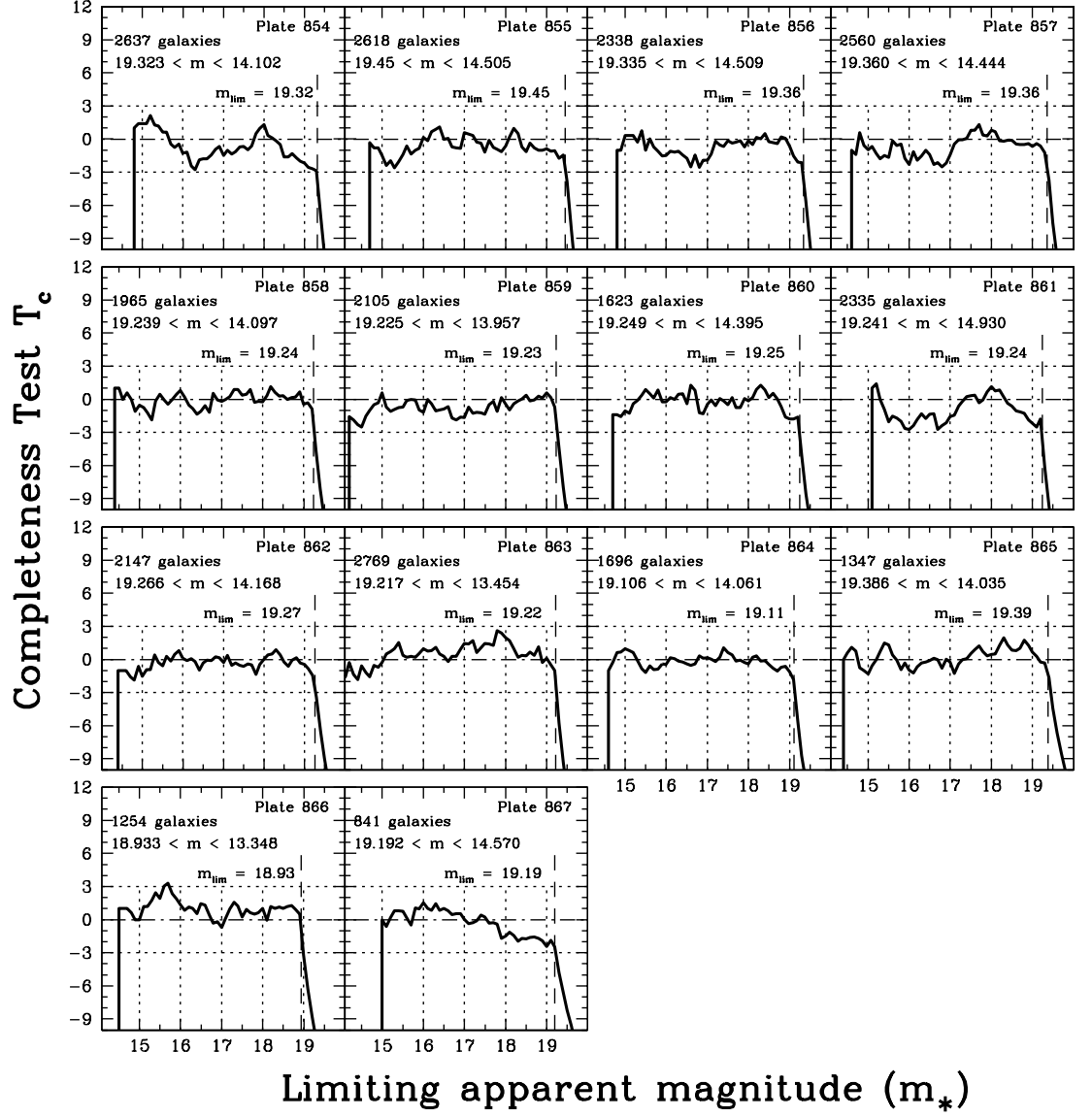
Figure 4.8: T_c curves for the 2dF-NGP plate numbers: 854 to 867.

Table 4.2: 2dF-Southern Plate Information. This table shows information regarding the Southern Galactic Pole region of 2dF. The *IFIELD* column refers to the UKST plate number. The following two columns show the number of galaxies in each plate before and after our selection criteria respectively. The faintest and brightest galaxies for each plate given in apparent magnitudes are those after selection.

Index	IFIELD	No. in IFIELD (raw)	No. in IFIELD (after selection)	Faintest app. magnitude	Brightest app. magnitude
1	349	3126	1722	19.431	14.441
2	350	2933	1524	19.455	14.508
3	351	1792	765	19.521	14.707
4	352	2288	913	19.430	14.387
5	353	2335	874	19.502	14.410
6	354	2936	1229	19.500	14.399
7	355	2635	1364	19.455	13.971
8	356	2094	948	19.519	14.553
9	357	1179	523	19.432	14.604
10	404	2465	1051	19.402	15.406
11	405	2843	1094	19.423	14.305
12	406	3296	1758	19.445	14.429
13	407	2783	1325	19.467	15.084
14	408	3260	1616	19.410	15.171
15	409	4576	2457	19.435	14.525
16	410	4269	2314	19.427	14.550
17	411	4429	2522	19.406	14.400
18	412	4378	1883	19.601	14.386
19	413	3956	1619	19.545	14.631
20	414	3840	1559	19.467	14.466
21	415	3998	2047	19.421	14.464
22	416	2798	1206	19.378	14.332
23	417	2512	1066	19.386	14.275
24	418	3932	1671	19.381	14.440
25	466	4054	2144	19.276	14.345
26	467	4268	2049	19.305	14.185
27	468	3186	1499	19.384	14.293
28	469	3928	2179	19.358	14.015
29	470	3821	1800	19.322	14.802
30	471	2671	1440	19.446	14.295
31	472	1816	879	19.440	14.691
32	473	1505	778	19.363	14.184
33	474	1698	1126	19.301	14.279
34	475	1361	576	19.365	14.706
35	476	1837	792	19.506	14.393
36	477	1415	510	19.342	14.601
37	478	1442	708	19.366	14.631
38	479	1710	762	19.400	13.893
39	480	2051	872	19.349	14.777
40	481	1880	1069	19.242	15.013
41	532	3663	1763	19.210	14.213
42	533	3319	1878	19.331	13.959
43	534	2454	1232	19.361	14.290
44	535	2575	1082	19.254	14.710
45	536	1945	873	19.307	14.845
46	537	933	528	19.390	14.885

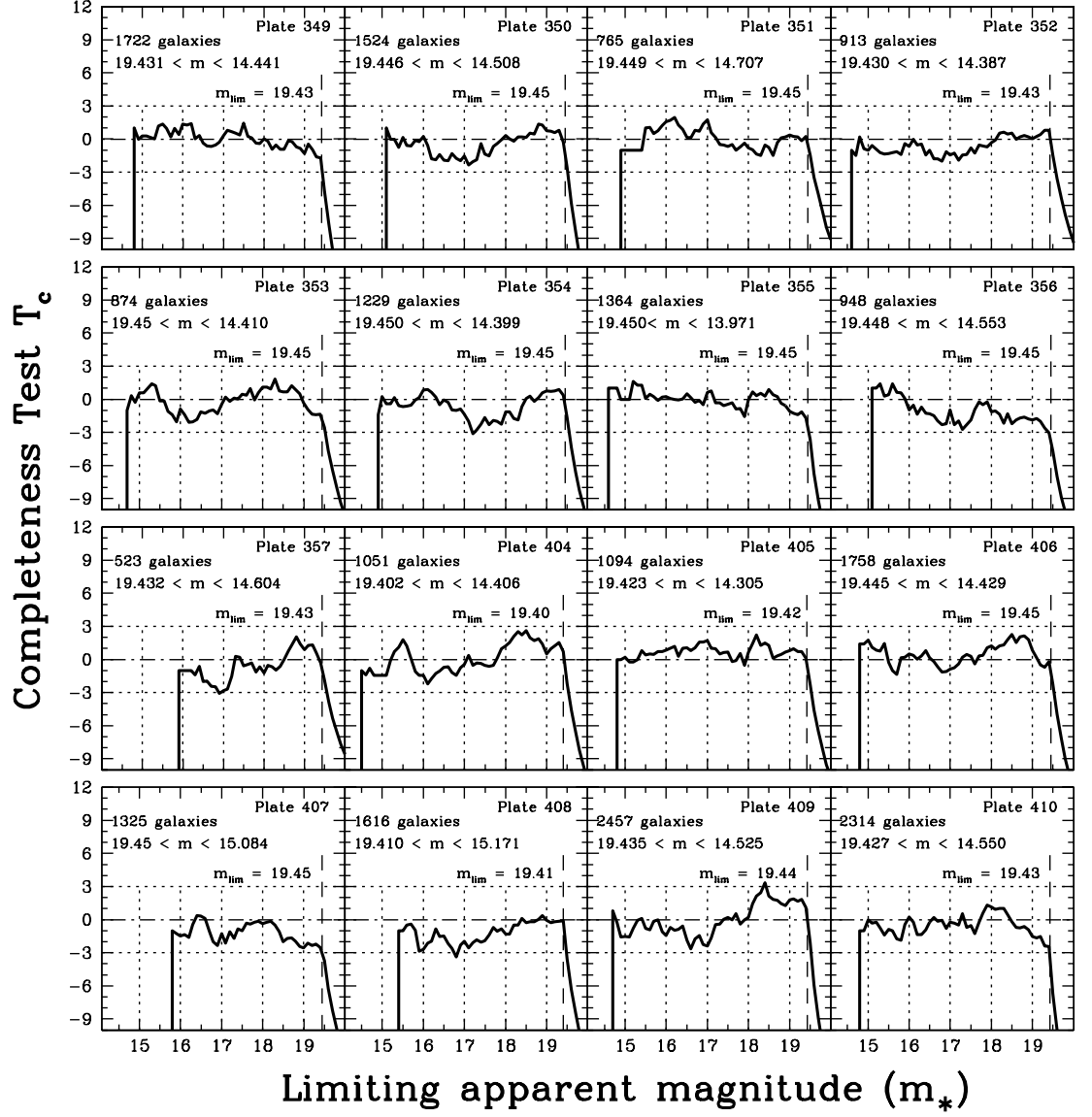


Figure 4.9: 2dF-SGP plates 349 to 410

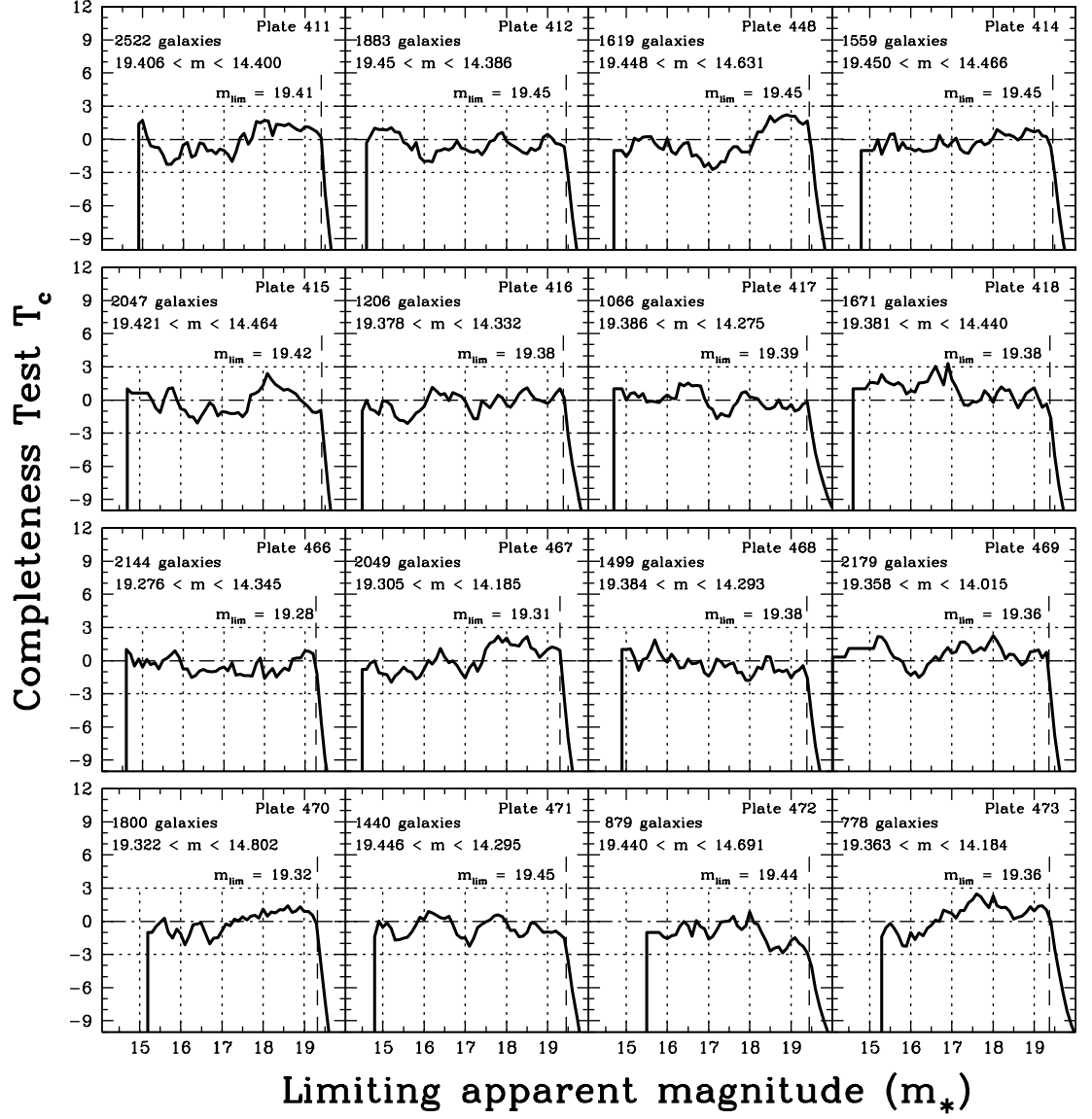


Figure 4.10: 2dF-SGP plates 411 to 473

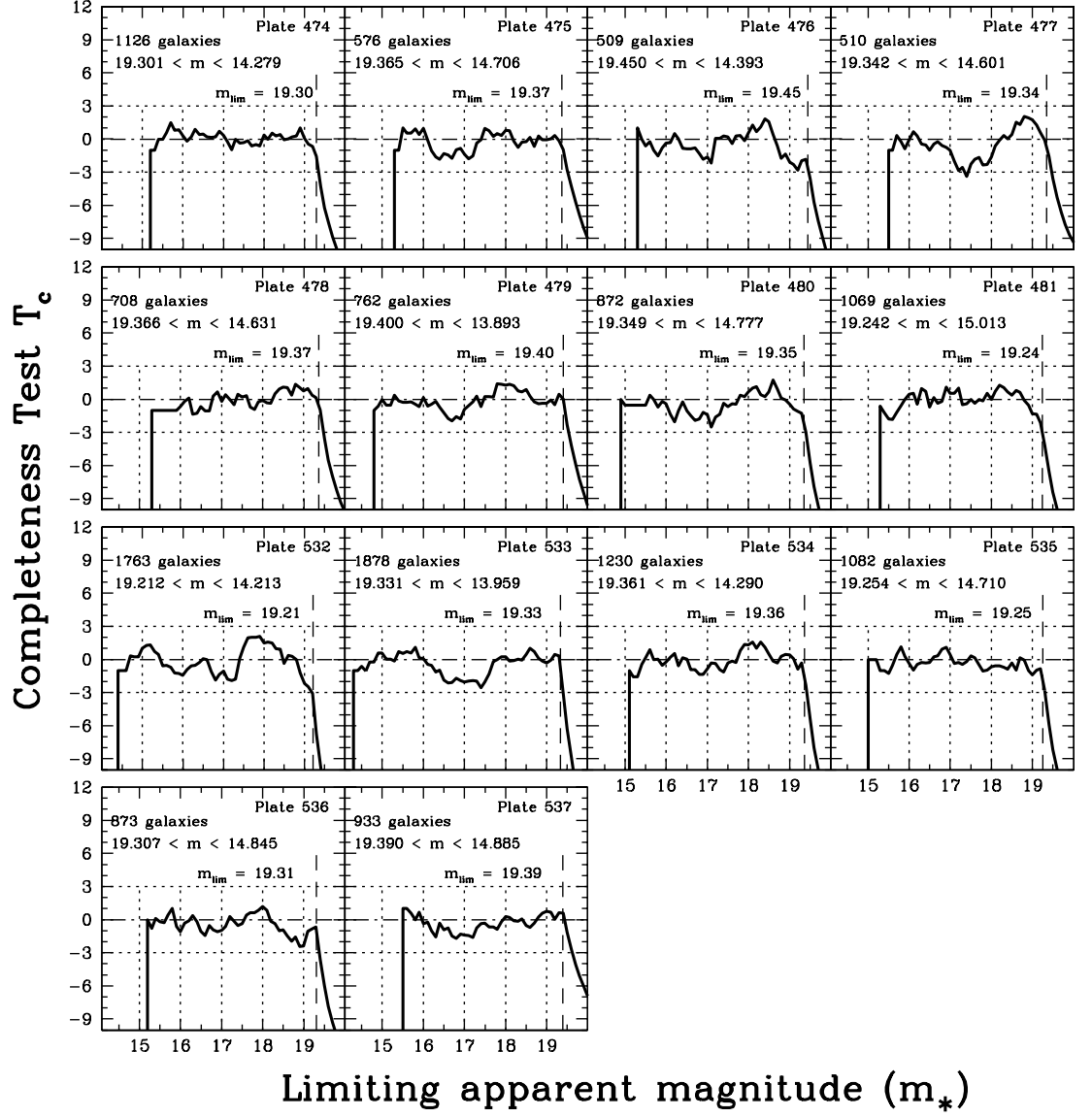


Figure 4.11: 2dF-SGP plates 474 to 537

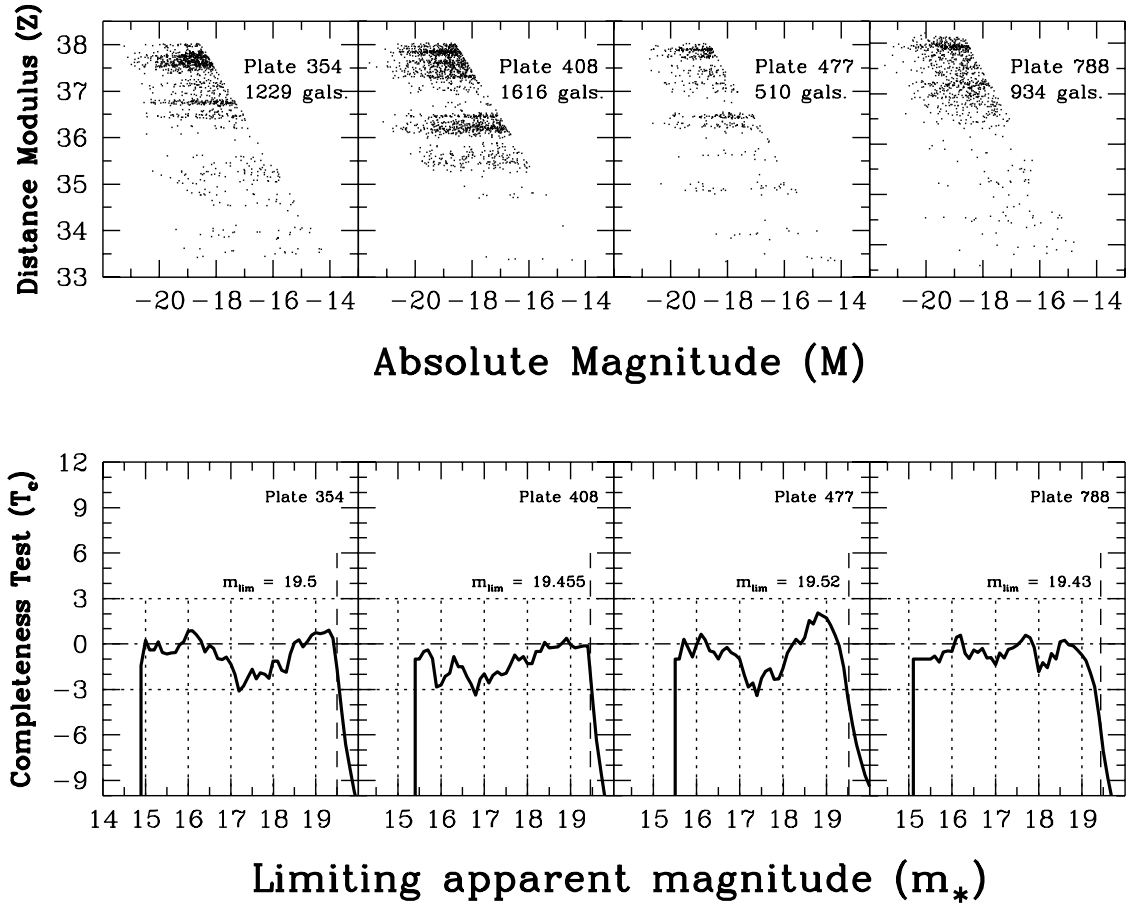


Figure 4.12: The first three plots, plates 354, 408 and 477 show the M - Z distribution (top) and the resulting T_c curve which all show a brief drop below the -3σ level. What is also clear is that all three curves show an overall similar trend to that of whole 2dF data-set in Figure 4.3. By examining the M - Z distributions of these plates we observe that all appear to have a well defined bright apparent magnitude. Plate 788 is randomly selected to show a T_c curve that behaves as one would expect for a complete sample. The M - Z distribution of plate 788 shows some scattering of nearby galaxies in the lower left of the figure, that would indicate there not being a clear bright limit.

4.2 Generalising the T_c statistic

4.2.1 Re-defining the random variable ζ

As in R01, the key element of our extended completeness test is the definition of a random variable, ζ , related to the cumulative luminosity function of the galaxy population. We proceed in a similar manner to R01, but now with both a bright and faint apparent magnitude limit.

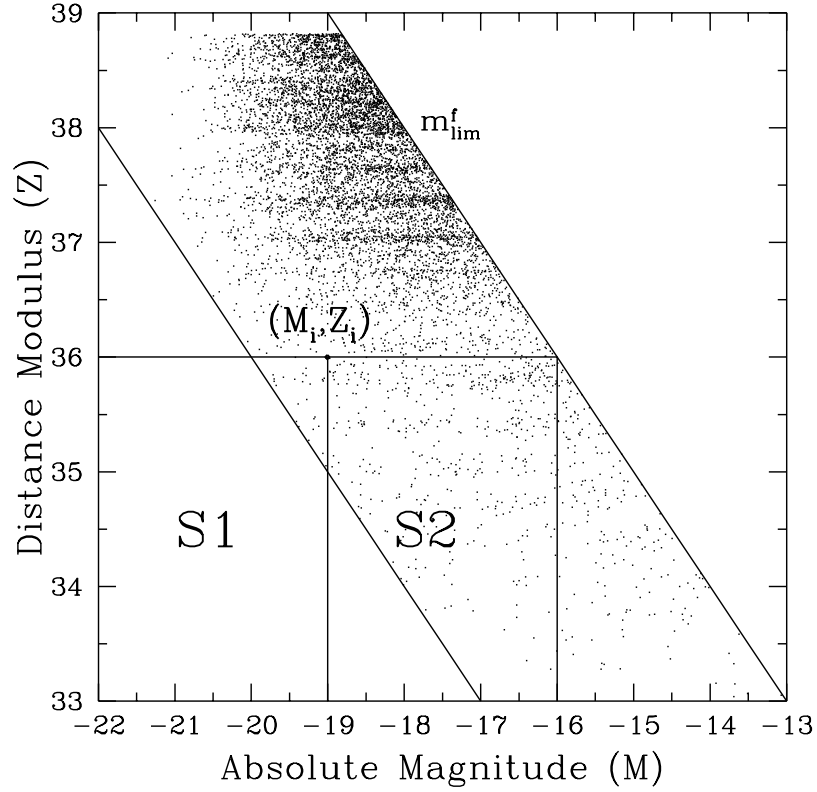


Figure 4.13: In this schematic we have imposed a bright magnitude limit of $m = 16.0$ mag to the MGC data-set to illustrate how the random variables M and Z remain unseparable under the traditional construct of S_1 and S_2 .

We now modify the selection function as previously defined by Equation 3.5 to include a bright limit in apparent magnitude, m_{lim}^b , such that

$$\psi(m, z, l, b) \equiv \theta(m_{\text{lim}}^f - m) \times \theta(m - m_{\text{lim}}^b) \times \phi(z, l, b), \quad (4.5)$$

Thus, the probability density function which was defined in equation 3.7 now takes form

$$dP = \bar{h}(Z)dZ f(M)dM \theta(m_{\text{lim}}^f - m)\theta(m - m_{\text{lim}}^b). \quad (4.6)$$

To see how we construct ζ in this more general case consider Figure 4.14, which schematically represents an M - Z plot of corrected distance modulus versus absolute magnitude for the observable population of galaxies. Shown in the plot are solid diagonal lines representing the ‘true’ faint and now the bright apparent magnitude limits, m_{lim}^f and m_{lim}^b respectively, while the red diagonal line represents putative faint magnitude limit, m_*^f . Comparing the original Rauzy construction of the completeness test from Figure 4.13, where a bright magnitude limit is well defined, with our new con-

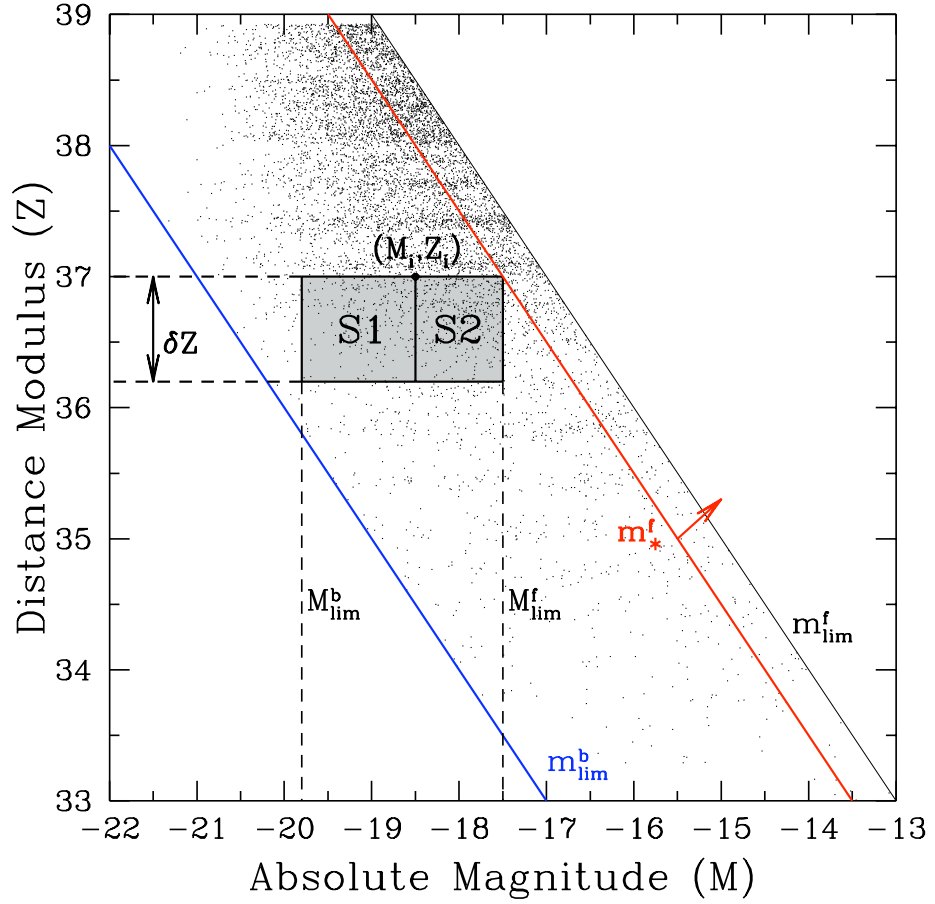


Figure 4.14: Diagram illustrating the construction of the rectangular regions S_1 and S_2 , defined for a typical galaxy at (M_i, Z_i) . The schematic shows the construction of regions S_1 and S_2 with the inclusion of a bright limit. These regions are uniquely defined for a slice of fixed width, δZ , in corrected distance modulus, and for ‘trial’ faint apparent magnitude limit m_*^f . Also shown are the true bright and faint apparent magnitude limits m_{lim}^b and m_{lim}^f , within which the rectangular regions S_1 and S_2 contain a joint distribution of M and Z that is separable.

struction in Figure 4.14, we can see that the addition of a bright magnitude limit has a major impact on the construction of this separable region: in short, the region is no longer unique.

In graphical terms, the essential idea of our extended completeness test is to identify from the data the faintest value of m_*^f and the brightest value of m_*^b which together bound a *rectangular* region of the M - Z plane, within which the joint distribution of M and Z for observable galaxies is separable. If we fix the width, δZ , in corrected distance modulus as shown in Figure 4.14, the corresponding separable region is now uniquely defined. Moreover we can then define for the i^{th} galaxy the following absolute

magnitudes:

- $M_{\text{lim}}^f(Z_i)$, the absolute magnitude of a galaxy, at corrected distance modulus Z_i , which would be observed at the true faint apparent magnitude limit m_{lim}^f ,
- $M_{\text{lim}}^b(Z_i - \delta Z)$, the absolute magnitude of a galaxy, at corrected distance modulus $Z_i - \delta Z$, which would be observed at the true bright apparent magnitude limit m_{lim}^b .

These two absolute magnitudes are indicated, for the corresponding *putative* faint magnitude limit m_*^f and the (assumed known) true bright magnitude limit, m_{lim}^b , by the vertical dashed lines in Figure 4.14.

We now re-define the random variable ζ as follows:

$$\zeta = \frac{F(M) - F[M_{\text{lim}}^b(Z - \delta Z)]}{F[M_{\text{lim}}^f(Z)] - F[M_{\text{lim}}^b(Z - \delta Z)]}, \quad (4.7)$$

where $F(M)$ is the CLF, i.e.

$$F(M) = \int_{-\infty}^M f(x) dx. \quad (4.8)$$

Thus ζ shares the same two defining properties as the corresponding random variable defined in R01. Equation 4.7 therefore generalises the definition of ζ to the case of a galaxy survey with bright and faint magnitude limits. The relevance of ζ as a diagnostic of magnitude completeness will be demonstrated in the next two sections.

4.2.2 Estimating ζ and computing the T_c statistic

As was the case in R01, the random variable ζ has the very useful property that we can estimate it without any prior knowledge of the CLF, $F(M)$. Given a value of δZ , it is clear from Figure 4.14 that for each point (M_i, Z_i) in the M - Z plane we can define the regions S_1 and S_2 as follows:

- $S_1 = \{(M, Z) : M_{\text{lim}}^b \leq M \leq M_i, Z_i - \delta Z \leq Z \leq Z_i\},$
- $S_2 = \{(M, Z) : M_i < M \leq M_{\text{lim}}^f, Z_i - \delta Z \leq Z \leq Z_i\}.$

In the special case where there is no bright limit the regions S_1 and S_2 are as shown in Figure 3.1.

Clearly the random variables M and Z are now independent within each sub-sample S_1 and S_2 . Therefore from Equation 4.6 the expected number of points r_i belonging to S_1 satisfies

$$\frac{r_i}{N_{\text{gal}}} = \int_{Z_i - \delta Z}^{Z_i} \bar{h}(Z') dZ' \times \int_{M_{\text{lim}}^b}^{M_i} f(M) dM, \quad (4.9)$$

where N_{gal} is the total number of galaxies in the sample. Similarly the expected number of points n_i in $S_i = S_1 \cup S_2$ satisfies

$$\frac{n_i}{N_{\text{gal}}} = \int_{Z_i - \delta Z}^{Z_i} \bar{h}(Z') dZ' \times \int_{M_{\text{lim}}^b}^{M_{\text{lim}}^f} f(M) dM. \quad (4.10)$$

The integrals over absolute magnitude in Equations 4.9 and 4.10 may be rewritten as

$$\int_{M_{\text{lim}}^b}^{M_i} f(M) dM = F[M_i(Z_i)] - F[M_{\text{lim}}^b(Z_i)], \quad (4.11)$$

and

$$\int_{M_{\text{lim}}^b}^{M_{\text{lim}}^f} f(M) dM = F[M_{\text{lim}}^f(Z_i)] - F[M_{\text{lim}}^b(Z_i)]. \quad (4.12)$$

Thus, given a pair of ‘trial’ magnitude limits m_*^f and m_*^b , it follows from Equation 4.7 and Equations 4.9 to 4.12 that an estimate of ζ for the i^{th} galaxy is simply the ratio of the number of galaxies belonging to S_1 and $S_1 \cup S_2$ respectively (where M_*^b and M_*^f replace M_{lim}^b and M_{lim}^f in the definition of S_1 and S_2). In fact an unbiased estimate of ζ for the i^{th} galaxy is (c.f. R01)

$$\hat{\zeta}_i = \frac{r_i}{n_i + 1}. \quad (4.13)$$

This estimator is identical to that defined in R01; the introduction of a bright magnitude limit has simply changed the definition of the random variable ζ itself and the membership criteria of the two regions S_1 and S_2 . Thus, provided that both $m_*^f \leq m_{\text{lim}}^f$ and $m_*^b \geq m_{\text{lim}}^b$, then under our null hypothesis $\hat{\zeta}_i$ will be uniformly distributed on $[0, 1]$ and uncorrelated with Z_i , exactly as was the case in R01. Moreover the expectation value E_i and the variance V_i of the $\hat{\zeta}_i$ are given respectively by

$$E_i = E(\hat{\zeta}_i) = \frac{1}{2}, \quad V_i = E\left[\left(\hat{\zeta}_i - E_i\right)^2\right] = \frac{1}{12} \frac{n_i - 1}{n_i + 1}. \quad (4.14)$$

Note that V_i tends towards the variance of a continuous uniform distribution between 0 and 1 when n_i is large.

As in R01, we can, therefore, combine the estimator $\hat{\zeta}_i$ for each observed galaxy into a single statistic, T_c , which we can use to test the magnitude completeness of our

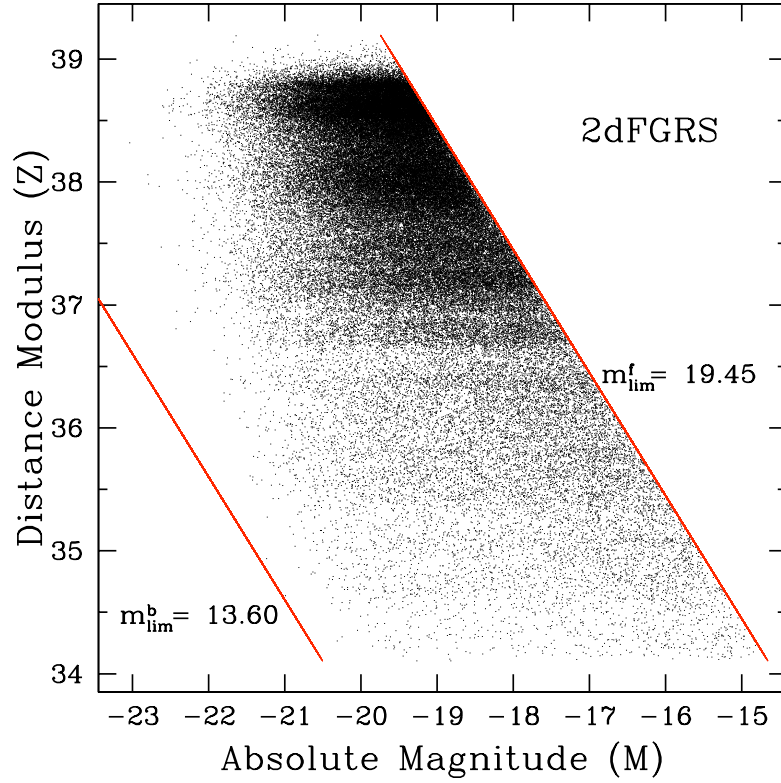


Figure 4.15: 2dFGRS M - Z distribution showing our adopted apparent magnitude bright limit of 13.60 mag.

sample for adopted trial magnitude limits m_*^f and m_*^b . T_c is thus defined as

$$T_c = \sum_{i=1}^{N_{\text{gal}}} \left(\hat{\zeta}_i - \frac{1}{2} \right) / \left(\sum_{i=1}^{N_{\text{gal}}} V_i \right)^{\frac{1}{2}}. \quad (4.15)$$

If the sample is complete in apparent magnitude, for a given pair of trial magnitude limits, then T_c should be normally distributed with mean zero and variance unity. If, on the other hand, the trial faint (bright) magnitude limit is fainter (brighter) than the true limit, T_c will become systematically negative, due to the systematic departure of the $\hat{\zeta}_i$ distribution from uniform on the interval $[0, 1]$.

4.3 Applying the Revised T_c to MGC and 2dFGRS

In the absence of a clear indication from the literature of what is an appropriate bright magnitude limit, we adopted the brightest galaxy in our main subset, $m_{\text{lim}}^b = 13.364$ mag. The right hand plot of Figure 4.16 shows the T_c curve obtained when we include

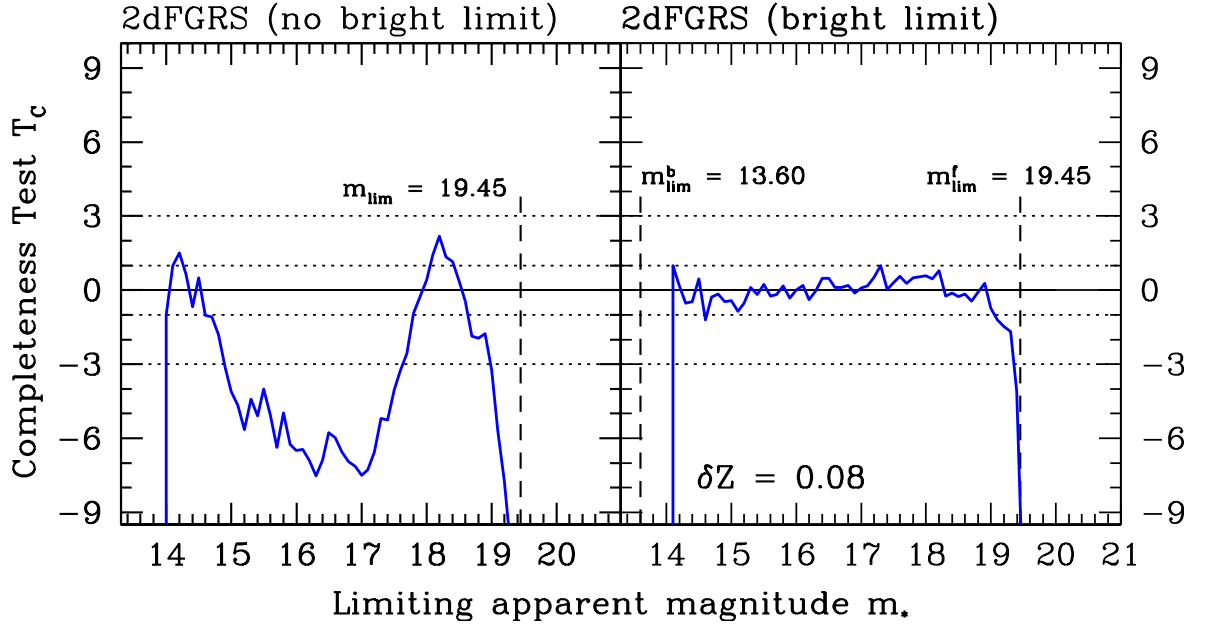


Figure 4.16: Performance of the T_c statistic applied to our 2dFGRS sample. In the left hand panel we compute T_c assuming only a faint magnitude limit, for both uncorrected magnitudes and after applying various different $(k + e)$ -corrections. In the right hand panel we include the effect of a bright apparent magnitude limit, adopting for simplicity a value equal to the apparent magnitude of the brightest galaxy in our sample. The resulting T_c curve (shown for uncorrected magnitudes and computed assuming a fixed ‘slice’ width of $\delta Z = 0.08$ in distance modulus) is now entirely consistent with magnitude completeness up to and including the published faint limit, but drops very sharply at fainter magnitudes.

this bright limit. The plot in this figure demonstrates a change in the behaviour of the 2dF completeness that is dramatic compared with the original R01 T_c statistic in Figure 4.16, left. Our choice of $\delta Z = 0.08$ was motivated by the fact that a small δZ leads to low numbers of galaxies within the subsets S_1 and S_2 , making our test statistic more sensitive to large statistical fluctuations and therefore less sensitive to a sharp cut in magnitude. We investigate the sensitivity to δZ more fully in Figure 4.17 and discuss the implications in greater detail in Chapter 6. However, we can see that by simply accounting for a bright limit - notwithstanding the fact that no published bright limit has been reported in the literature - we find that the 2dFGRS data-set is indeed complete to the published faint magnitude limit with no evidence for residual systematics. Having already established in Chapter 2 that MGC is indeed complete up to the published faint magnitude limit of 20.0 mag, we can now use this survey to demonstrate how the presence of a well defined bright limit can affect the Rauzy completeness test, if not properly accounted for as was described in Section 4.1.4.

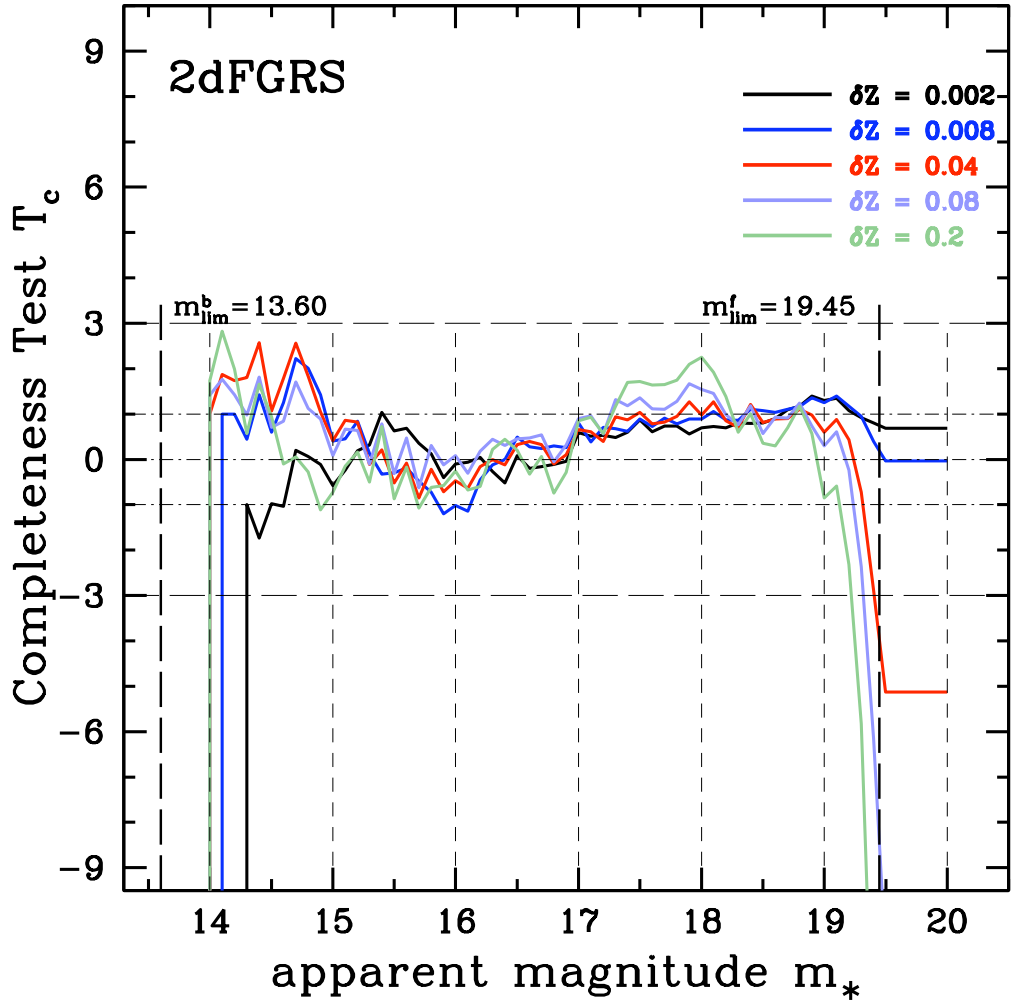


Figure 4.17: Our generalised version of T_c applied to 2dFGRS showing the effect of varying values of δZ . As can be seen, for very small values of δZ i.e 0.002 and 0.008 in this case, the statistic suffers from small numbers with the effect of T_c ‘flat-lining’ within the 3σ level. As we increase the value of δZ the T_c statistic begins to drop sharply for values of $m_* > m_{lim}^f$.

To do this we imposed: $m_{lim}^b > 14$, $m_{lim}^b > 15$, and $m_{lim}^b > 16$ mags respectively (see Figure 4.18). Figure 4.19 - left shows T_c curves for the data-sets with these artificial bright limits, but where T_c was computed assuming *no* bright limit. The plots clearly show that, as the bright limit is made progressively fainter, the computed value of T_c deviates more strongly from the behaviour expected for a complete data-set. This trend is as expected: as can be seen in Figure 4.19 - left, the presence of the bright limit breaks the separability of the M and Z distributions for observable galaxies. Hence, if the bright limit is ignored then the computed value of T_c will be systematically biased.

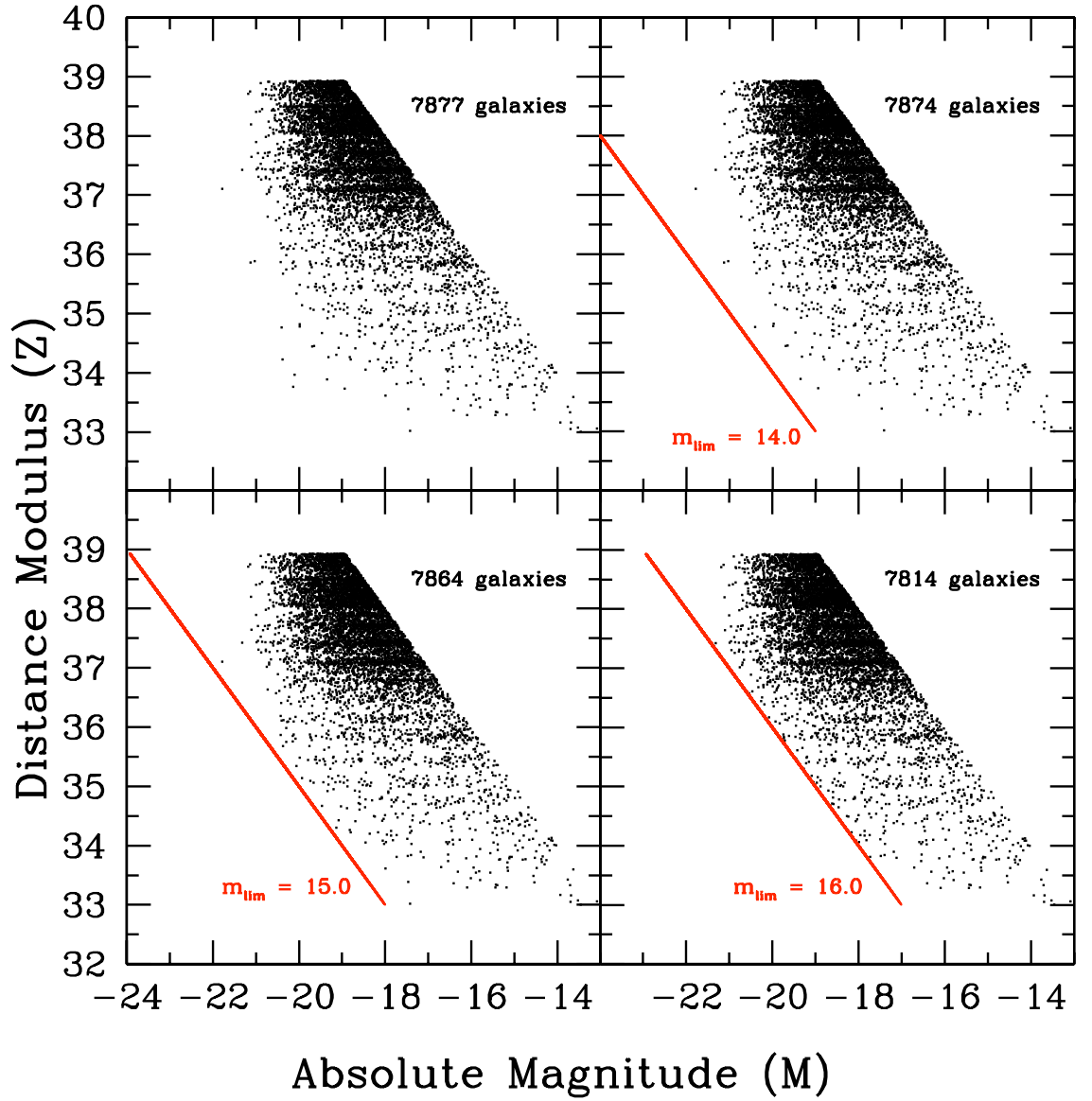


Figure 4.18: MGC M - Z distributions with progressively fainter bright limits: top left - no cut in apparent magnitude, top right - $m > 14$, bottom left - $m > 15$ and bottom right - $m > 16$ respectively.

We now impose the same artificial bright limits as before but apply our generalised T_c method which accounts for a bright *and* faint limit (see Figure 4.19 - right) and fix the value for m_*^b in each slice to be equal to the brightest observed galaxy. It is evident from this plot that, even with a bright magnitude limit as faint as $m = 16$ mag, the performance of the T_c statistic at fainter magnitudes is largely unaffected, showing consistent behaviour for *all* the bright limits considered.

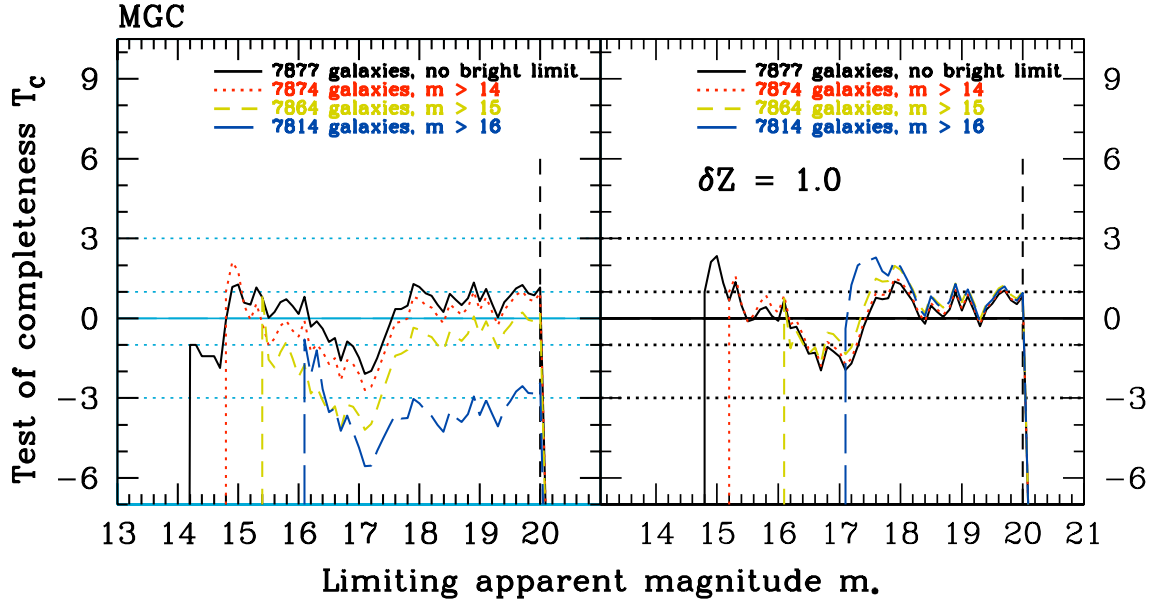


Figure 4.19: The T_c statistic computed for the MGC survey (without $(k+e)$ -corrections) but now illustrating the effect, close to the faint limit, of imposing artificially a bright apparent magnitude limit on the selected galaxies. In the left panel the solid black curve shows T_c computed assuming no bright limit (identical to the right hand portion of the red dashed curve in Figure 3.5 on page 62 while the other three curves correspond to progressively fainter bright limits, at $m > 14$, $m > 15$ and $m > 16$ mag respectively). For all four curves we calculated T_c following R01 – i.e. assuming *no* bright limit. We can clearly see that the presence of a bright limit, if ignored, has a significant impact on the computed value of T_c for faint magnitudes, and thus could adversely affect the assessment of magnitude completeness close to the faint limit. In the right hand panel we repeat our analysis for the same four cases as in the left panel, but now use our extended method which explicitly accounts for the presence of the bright limit. We can clearly see that the performance of T_c is no longer adversely affected, and a consistent estimate of the faint magnitude limit is obtained for different imposed bright limits.

4.4 Analysis of the Sloan Digital Sky Survey - Early Types

4.4.1 The data

The Sloan Digital Sky Survey (SDSS) used the 2.5m Ritchey-Chretien wide-field altitude-azimuth telescope located at the Apache Point Observatory in New Mexico. For the spectroscopy the SDSS team utilise a pair of spectrographs capable of measuring more than 600 redshifts simultaneously. See [Stoughton \(2002\)](#) for a description of the Early Data Release; [Abazajian \(2003\)](#) for a description of DR1, the First Data Release; [Gunn \(1998\)](#) for a detailed description of the camera; [Fukugita et al. \(1996\)](#), [Hogg](#)

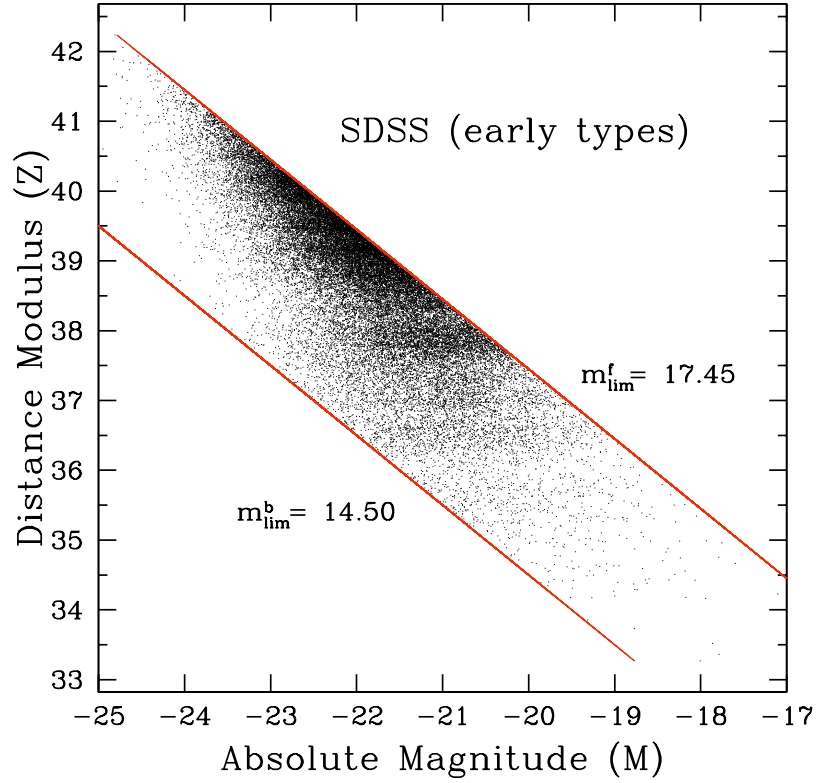


Figure 4.20: The M - Z distribution for the SDSS-Early Types. This particular data set has well defined apparent magnitude limits as indicated by the red lines. To correctly apply the T_c statistic, the traditional Rauzy test had to be extended to include the bright limit in apparent magnitude.

et al. (2001) and Smith (2002) for details of the photometric system and calibration; Lupton (2001) for a discussion of the photometric data reduction pipeline.

There has been a total of seven data releases since 2001. See York (2000) for a technical summary of the SDSS project; Pier et al. (2003) for the astrometric calibrations; Blanton et al. (2003) for details of the tiling algorithm; Strauss (2002) and Eisenstein (2001) for details of the target selection.

In broad terms, the SDSS sample includes spectroscopic information as well as photometric measurements in the u^* , g^* , r^* , i^* and z^* bands. The SDSS First Data Release covers an area of $\approx 2000 \text{ deg}^2$ (Abazajian, 2003) on the sky.

4.4.2 Selection limits and cosmology

For our analysis we used galaxies present in the SDSS, early types only (hereinafter referred to as ‘SDSS-Early Types’). The selection criteria have been discussed in detail

in Bernardi (2003a) and the data-set compiled in Bernardi et al. (2005). A total of 39,320 objects have been targeted as early-type galaxies and have dereddened Petrosian (hereinafter referred to as, m_r) apparent magnitudes in the range $14.5 < m_r < 17.75$, and a redshift range of $0.0 < z < 0.4$.

We have assumed a Hubble constant of $70 \text{ kms}^{-1} \text{ Mpc}^{-1}$ as adopted in Bernardi (2003b). It should be noted that within the range of H_0 values that we have considered throughout this thesis i.e. 70 to $100 \text{ kms}^{-1} \text{ Mpc}^{-1}$, the effect on our statistics has been minimal.

4.4.3 Results

As previously discussed, the SDSS data-set has both a published bright and faint apparent magnitude limit of $m_r = 14.55$ and $m_r = 17.45$ mag respectively. Figure 4.20 shows the corresponding $M-Z$ distribution for the selection we are considering. We, therefore, tested the completeness of the SDSS-Early Type galaxies using our generalised T_c statistic which accounts for both a faint and bright limit. As an illustration, we chose to fix the bright limit of the SDSS data-set to be equal to the published value, and computed the T_c statistic as a function of the trial faint magnitude limit. We set $\delta Z = 0.2$. Figure 4.21-left shows the resulting T_c curve. We see that our results are in agreement with the published faint magnitude limit – i.e. the behaviour of T_c is consistent with magnitude completeness up to and including a sharp, faint limit of $m_r = 17.45$ mag, followed immediately by the strongly negative behaviour expected for an incomplete sample at fainter magnitudes.

The right-hand panel of Figure 4.21 shows the T_c curve computed using the traditional Rauzy method with a faint limit only. Here the results show a similar trend to that seen for the MGC data-set with an artificially imposed bright limit section, and that of the 2dFGRS before a bright apparent magnitude limit was adopted. The results therefore further underlines the importance of correctly accounting for a bright *and* faint limit when both are present in the data.

4.5 Conclusions

In this chapter we have shown that our initial approach to the 2dFGRS was to apply the original Rauzy (2001) test which accounts for a single, faint magnitude limit only. This approach was motivated by the current literature, in which *only* a faint limiting magnitude of $m_{\text{lim}}^f = 19.45$ mag was defined for the survey. However, the application

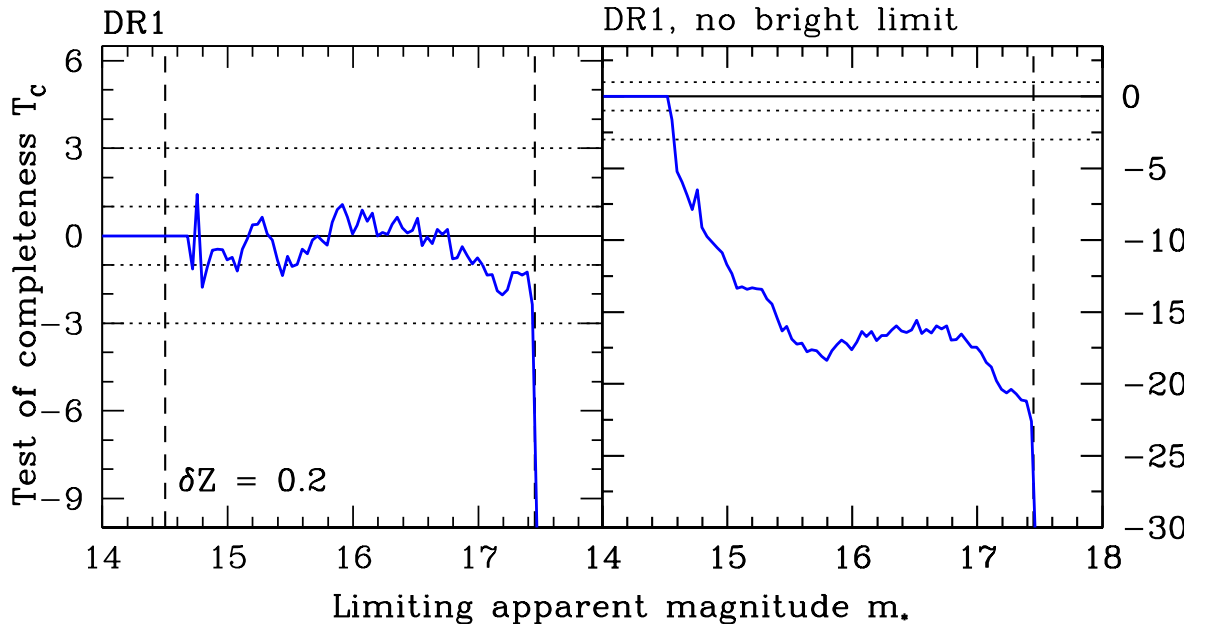


Figure 4.21: Performance of the T_c statistic applied, as an illustrative example, to the SDSS-Early Type elliptical galaxies. In the left panel we compute T_c using our extended method, fixing the bright magnitude limit to equal the published value of 14.55 mag. We see that, in this case, the behaviour of T_c is consistent with magnitude completeness up to and including the published faint limit of 17.45 mag., but the statistic drops rapidly thereafter – indicating the sharp onset of magnitude incompleteness. In the right hand panel, on the other hand, we compute T_c following Rauzy (2001) – i.e. assuming a faint magnitude limit only. As anticipated, we see that the test statistic deviates very strongly from its expected value for a complete data-set at magnitudes which are much brighter than the published faint magnitude limit (although it is worth noting that T_c still decreases even more rapidly once the published faint limit is exceeded).

of our T_c test revealed that the 2dF survey is strongly inconsistent with being complete in apparent magnitude when only a sharp, faint limit is adopted. Furthermore, the SDSS-Early Types data-set we used has well defined bright and faint apparent magnitude limits and therefore we developed the completeness test statistic, T_c , technique to account for the presence of both a faint and bright apparent magnitude limit in magnitude-redshift samples.

Applying our generalised method to SDSS-Early Types, 2dFGRS and MGC surveys confirms the completeness of data-sets such as SDSS-Early Types where a faint and bright limit is well defined and published in the literature. Specifically, we have demonstrated that the SDSS-Early Types are complete in apparent magnitude up to the published magnitude limit of $m_r = 17.45$ mag indicating no residual systematics.

Similarly, we confirm the 2dFGRS indicated completeness up to the $m_{\text{lim}}^f = 19.45$ mag

but *only* if a secondary bright limit ($m_{\text{lim}}^{\text{b}} = 13.6$ mag in our subset) is also included.

Lastly, using the 2dFGRS data, we observed that for small values δZ the T_c statistic will be dominated by shot-noise, resulting in a flat-line effect (or seemingly *over* completeness) beyond the magnitude limit. We shall explore this behaviour in greater detail in Chapter 7.

Chapter 5

Introducing the T_v statistic

Rimmer - *“Step up to red alert.”*

Kryten - *“Sir, are you absolutely sure? It does mean changing the bulb.”*

From the series *Red Dwarf*

In this brief chapter we introduce a further variant on the test statistic T_c , which we have named T_v , that is related to the distribution of corrected distance modulus for observable galaxies in a magnitude-redshift survey and which is actually close in spirit to the V/V_{\max} test.

5.1 Construction the T_v statistic

5.1.1 Defining the random variable τ

Figure 5.1 shows two schematic M - Z plots that are analogous to Figures 3.1 and 4.14. Therefore, we consider the following two cases:

- **Case I:** Data-sets with a single faint apparent magnitude limit, m_{lim}^f , in keeping with R01,
- **Case II:** Datasets with both a faint *and* bright limit, m_{lim}^f and m_{lim}^b respectively.

The left hand panel in Figure 5.1 illustrates **Case I** with the ‘putative’ faint limit, m_*^f shown as a red diagonal line and as before. The right hand plot illustrates **Case II** showing a ‘true’ bright and faint apparent magnitude limit, m_{lim}^b (blue diagonal line) and m_{lim}^f , with ‘putative’ faint limit, m_*^f , shown as the bold red diagonal line. Again, the position, (M_i, Z_i) , of a typical galaxy is shown on each panel.

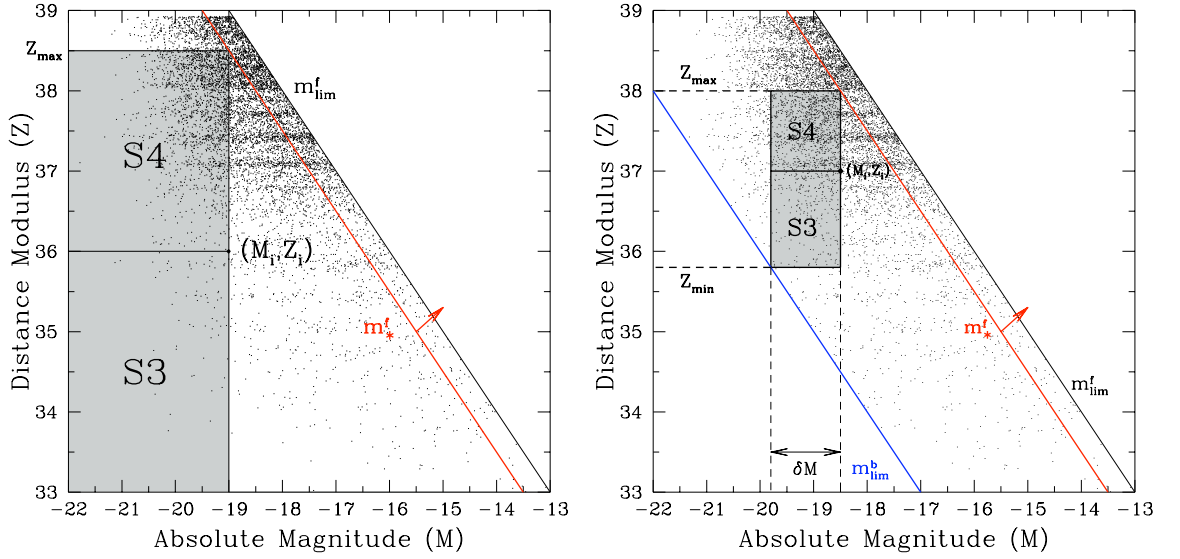


Figure 5.1: Schematic diagram illustrating the construction of the rectangular regions S_3 and S_4 , defined for a typical galaxy at (M_i, Z_i) , which feature in the estimation of our new completeness test statistic, T_v . The left hand panel illustrates how S_3 and S_4 are constructed for a survey with only a faint magnitude limit m_{lim}^f , and are shown for a trial faint limit m_*^f . The right hand panel shows the case where the survey also has a true bright limit m_{lim}^b (which we assume for simplicity is known), and the rectangles are constructed for trial bright and faint limits m_{lim}^b and m_{lim}^f respectively. Note that the construction of the rectangles is unique for a ‘slice’ of fixed width, δM , in absolute magnitude.

We can define for the i^{th} galaxy:

- **Case I and II** : $Z_{\text{max}}(M_i)$, the corrected distance modulus of a galaxy, with absolute magnitude M_i , which would be observed at the true faint apparent magnitude limit m_{lim}^f .

In the right hand panel of Figure 5.1, we consider a ‘slice’ of width δM in absolute magnitude brighter than M_i . Therefore, for **Case II** we see that, given the three quantities, m_{lim}^b , m_{lim}^f and δM , we can also define the following corrected distance moduli such that:

- **Case II only** : $Z_{\text{min}}(M_i - \delta M)$, the corrected distance modulus of a galaxy, with absolute magnitude $M_i - \delta M$, which would be observed at the true bright apparent magnitude limit m_{lim}^b .

These two limiting distance moduli are indicated, for the *putative* faint magnitude limit m_*^f and the true bright magnitude limit, m_{lim}^b , by the horizontal dashed lines in Figure 5.1.

We now let $H(Z)$ denote the cumulative distribution function of corrected distance modulus for observable galaxies, i.e.

$$H(Z) = \int_{-\infty}^Z \bar{h}(Z') dZ'. \quad (5.1)$$

Then τ is defined as

$$\text{Case I : } \tau = \frac{H(Z)}{H[Z_{\max}(M)]}. \quad (5.2)$$

$$\text{Case II : } \tau = \frac{H(Z) - H[Z_{\min}(M - \delta M)]}{H[Z_{\max}(M)] - H[Z_{\min}(M - \delta M)]}. \quad (5.3)$$

It is straightforward to show that in both these cases τ possesses the following properties:

- P1: τ is uniformly distributed between 0 and 1,
- P2: τ and M are statistically independent.

These two properties are exactly analogous to the defining properties of ζ , except that τ is now independent of the distribution of corrected absolute magnitude, M . Once again we can use property P1 to construct a test for completeness in apparent magnitude.

5.1.2 Estimating τ and computing the T_v statistic

Under the assumptions introduced in the previous section, it follows that τ can be estimated from our observed data without any prior knowledge of the spatial distribution of galaxies. To see how this estimate is constructed, consider again Figure 5.1. For each point with co-ordinates (M_i, Z_i) in the M - Z plane we can define the regions S_3 and S_4 as follows:

Case I :

- (i) $S_3 = \{(M, Z) : M \leq M_i, Z \leq Z_i\}$,
- (ii) $S_4 = \{(M, Z) : M \leq M_i, Z_i \leq Z \leq Z_{\max}^i\}$.

Case II :

- (i) $S_3 = \{(M, Z) : M_i - \delta M \leq M \leq M_i, Z_{\min} \leq Z \leq Z_i\}$,
- (ii) $S_4 = \{(M, Z) : M_i - \delta M \leq M \leq M_i, Z_i \leq Z \leq Z_{\max}^i\}$.

It should be noted that in **Case I**, S_3 is identical to the region S_1 shown in Figure 3.1.

As in Chapter 2, we see that the random variables M and Z are independent in each sub-sample S_3 and S_4 . Therefore we can estimate τ by counting the number, r_i , of galaxies that belong to S_3 and the number, t_i , of galaxies that belong to $S_3 \cup S_4$. Similarly, an unbiased estimate of τ is given by

$$\hat{\tau}_i = \frac{r_i}{t_i + 1}. \quad (5.4)$$

Thus, provided that $m_*^f \leq m_{\text{lim}}^f$, then under our null hypothesis $\hat{\tau}_i$ will be uniformly distributed on $[0, 1]$ and uncorrelated with M_i , exactly analogous to the properties of for $\hat{\zeta}_i$ under the same hypothesis. Moreover the expectation E_i and variance V_i of the $\hat{\tau}_i$ are respectively

$$E_i = E(\hat{\tau}_i) = \frac{1}{2}, \quad V_i = E[(\hat{\tau}_i - E_i)^2] = \frac{1}{12} \frac{t_i - 1}{t_i + 1}. \quad (5.5)$$

Again, the variance of $\hat{\tau}_i$ tends towards that of a continuous uniform distribution between 0 and 1 for large t_i .

We can, therefore, again combine the estimator $\hat{\tau}_i$ for each observed galaxy into a single statistic, T_v , which we can use to test the magnitude completeness of our sample for adopted trial magnitude limits m_*^f and m_*^b . T_c is defined as

$$T_v = \sum_{i=1}^{N_{\text{gal}}} \left(\hat{\tau}_i - \frac{1}{2} \right) / \left(\sum_{i=1}^{N_{\text{gal}}} V_i \right)^{\frac{1}{2}}. \quad (5.6)$$

If the sample is complete in apparent magnitude, for a given pair of trial magnitude limits, then T_v should be normally distributed with mean zero and variance unity. If, on the other hand, the trial faint magnitude limit is fainter than the true limit, in either case T_v will become systematically negative, due to the systematic departure of the $\hat{\tau}_i$ distribution from uniform on the interval $[0, 1]$.

5.2 Application of the T_v Statistic

In the previous chapters we introduced and applied our improved T_c statistic, which can account for both a faint and bright magnitude limit in assessing the completeness of a magnitude-redshift survey. In this section we apply the T_v statistic, introduced in § 5.1 above, to the same data-sets. Our T_v statistic can be thought of as an improved, differential, version of the classical V/V_{max} test of galaxy evolution, which is generally presented in the literature as yielding a single number – the mean value of V/V_{max}

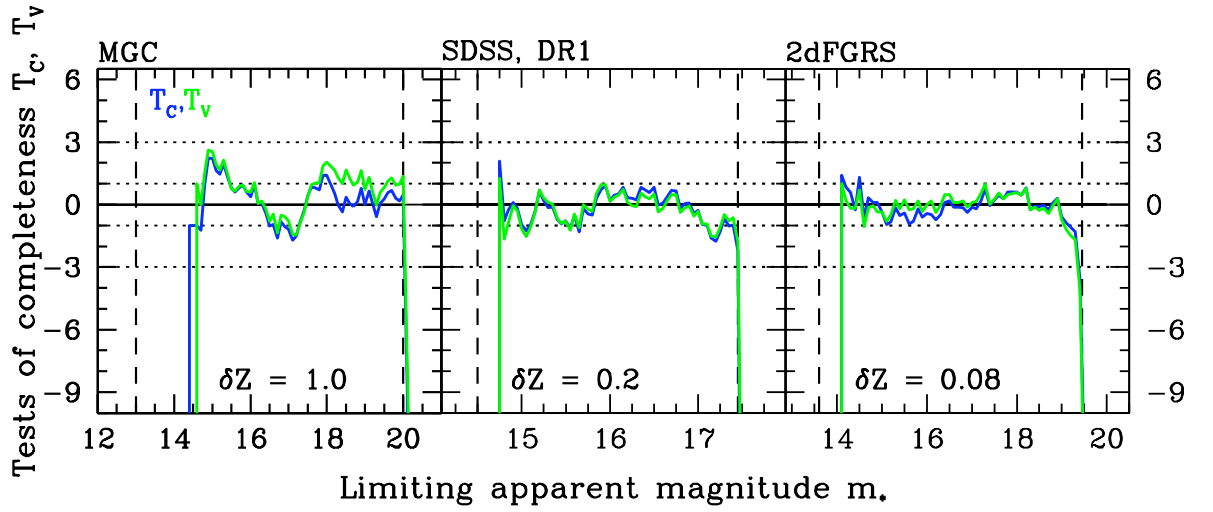


Figure 5.2: Comparison of the T_c and T_v statistics computed for MGC (left panel), SDSS-Early Types (middle panel) and 2dFGRS (right panel). For the latter two surveys the same bright limits were adopted as in Figures 4.21 and 4.16 respectively, and appropriate values for δZ (for T_c) and δM (for T_v) were chosen. Note the almost identical agreement of the test statistics in each case. To illustrate the robustness of this result, the MGC results are with $(k+e)$ -corrections applied, the SDSS-Early Types results are with k -corrections only applied, while the 2dFGRS results are for uncorrected galaxy data.

averaged over all galaxies in the survey, adopting a given faint apparent magnitude limit and assuming that the underlying spatial distribution of galaxies is homogeneous. In contrast, we can compute T_v as a function of an incrementally increasing m_\star^f and thus analyse our data-set via a series of progressively truncated subsets.

Figure 5.2 shows a comparison of the T_v and T_c curves for all three surveys. The left hand plot is the MGC survey with $(k+e)$ -corrections applied. The T_v curve shows an almost identical match to the T_c statistic. Similar behaviour is evident with SDSS-Early Types and 2dFGRS (middle and right plots). That T_v and T_c give a consistent indication of the completeness of these surveys should not be too surprising, since we are confident (at least once a bright limit is included in our analysis of the 2dFGRS) that all three are well calibrated and well understood. Moreover, they are all relatively shallow in redshift range, which means that extinction and evolution corrections are not likely to impact too strongly on our assessment of their completeness (a fact which is supported by our results for T_c). However, one can ask under what conditions might the two statistics T_c and T_v diverge from each other?

Consider a galaxy, i , characterised by its ‘coordinates’ (M_i, Z_i) . We have two complementary ways of generating magnitude limited data-sets for such a pair:

- At fixed luminosity we can ask what redshift distribution will produce apparent

magnitudes permitted by our selection criteria?

- Alternatively, at fixed redshift we can ask what distribution of luminosities (i.e. what part of the underlying galaxy luminosity function) are we sampling, given our selection limits in apparent magnitude?

The former criterion resembles the procedure used to construct the T_v statistic, while the latter criterion is more closely related to the procedure used to construct T_c . This also implies that one might expect the two statistics to behave differently when evolution becomes important – simply because evolution will, of course, break the separability of the *underlying* joint distribution of M and Z , i.e. the conditional distributions of M at given Z and Z at given M will no longer be simply equal to their marginal distributions. It seems likely, therefore, that an exploration of the systematic differences between T_c and T_v for deeper surveys may be an effective probe of evolution. We will investigate this further in subsequent chapters.

5.3 Conclusion

We have demonstrated in this chapter the development of a variant on the original Rauzy T_c completeness statistic, which we denote by T_v , based on the cumulative distance distribution of galaxies in a magnitude-redshift survey. We find that T_v has potential advantages over the widely used V/V_{\max} test: not least, the T_v statistic retains the same properties as that of T_c – i.e. is independent of the spatial distribution of galaxies within the survey. Furthermore, we have shown by example, that T_v – when applied to the same well calibrated and relatively shallow survey samples as T_c – produces almost identical results to that of the T_c statistic.

Chapter 6

Analysis of the CCLQG Survey: GALEX Selected Sample

“I can prove anything by statistics except the truth.”

George Canning - From *A Dictionary of Thoughts* (1908)

6.1 The Data and Sample Selection

This data-set has been compiled from two slightly overlapping 1.2 degree fields within the Clowes-Campusano Large Quasar Group (herein referred to as CCLQG) in the Far-UV (FUV; $\lambda_{eff} = 1538.6\text{\AA}$) and Near-UV (NUV; $\lambda_{eff} = 2315.7\text{\AA}$) filter bands, using the UV satellite GALEX (GALaxy Evolution eXplorer). The pipeline reduction was done by the GALEX team, including the photometric calibration. Optical complementary data are from the Sloan Digital Sky Survey DR5 ([Adelman-McCarthy and *et al.*, 2007](#)), which is sensitive to limiting apparent magnitudes of $u = 22.0$, $g = 22.2$, $r = 22.2$, $i = 21.3$, $z = 20.5$.

The source catalogue consists of 15688 sources created using SExtractor v.2.5.0 ([Bertin and Arnouts, 1996](#)). By excluding bad regions at the edge of the GALEX images and saturated sources using weight maps (WEIGHT WATCHER VERSION 1.7) as well as flag images, a cross-correlation with the SDSS-DR5 resulted in 14316 sources (matching radius: 4.5 arcsec).

The data was further reduced to clean the sample from false detections (e.g. bright star contaminations and reflections) such that only objects which have a SExtractor

extraction $FLAGS \leq 2$ in the NUV filter were selected, which resulted in a subsample of 13760 objects (final UV selected sample). Finally, we consider all galaxies in the range $0.05 < z_{\text{photo}} < 2.5$ which leaves us with a final sample of 13748 objects in which the brightest and faintest galaxies are $m^b = 14.33$ and $m^f = 25.72$ respectively in the NUV filter.

The raw apparent magnitudes have already been corrected for galactic extinction, $A_g(l, b)$, and we therefore convert to absolute magnitudes by,

$$M_i = m_i - 5 \log_{10}(d_{L_i}) - 25 \quad (6.1)$$

where there are as yet no published evolution correction. The distance modulus, Z is calculated by,

$$Z_i = 5 \log_{10}(d_{L_i}) + 25 \quad (6.2)$$

where the luminosity distance, d_{L_i} , is given by Equation 3.2 on page 54. We also adopt a Hubble constant of $H_0 = 100 \text{ kms}^{-1} \text{ Mpc}^{-1}$. Figure 6.1 shows the M - Z distribution of our final survey sample.

6.2 Results

The completeness test results for this selected data-set revealed two very interesting features shown in Figure 6.2:

1. a large spike in the T_v curve and,
2. a systematic drop in both T_c and T_v , below -3σ considerably brighter than the magnitude limit of the survey.

We shall explore both features separately to allow an accurate assessment of completeness of the data.

6.2.1 Photometric redshift truncation effect

The T_c and T_v results in Figure 6.2 show the original R01 method i.e. assuming a faint limit only. What is immediately clear is an uncharacteristic departure of T_v (the red curve) from T_c (the black curve) that begins at $m_* \approx 18.7 \text{ mag}$. The T_v statistic then rises sharply peaking at $m_* = 22.6 \text{ mag}$, $\sigma = 18.58$ before dropping systematically below -3σ at $m_* \approx 24.45 \text{ mag}$. A result with this atypical behaviour has not been observed for any of the surveys that we have examined thus far and

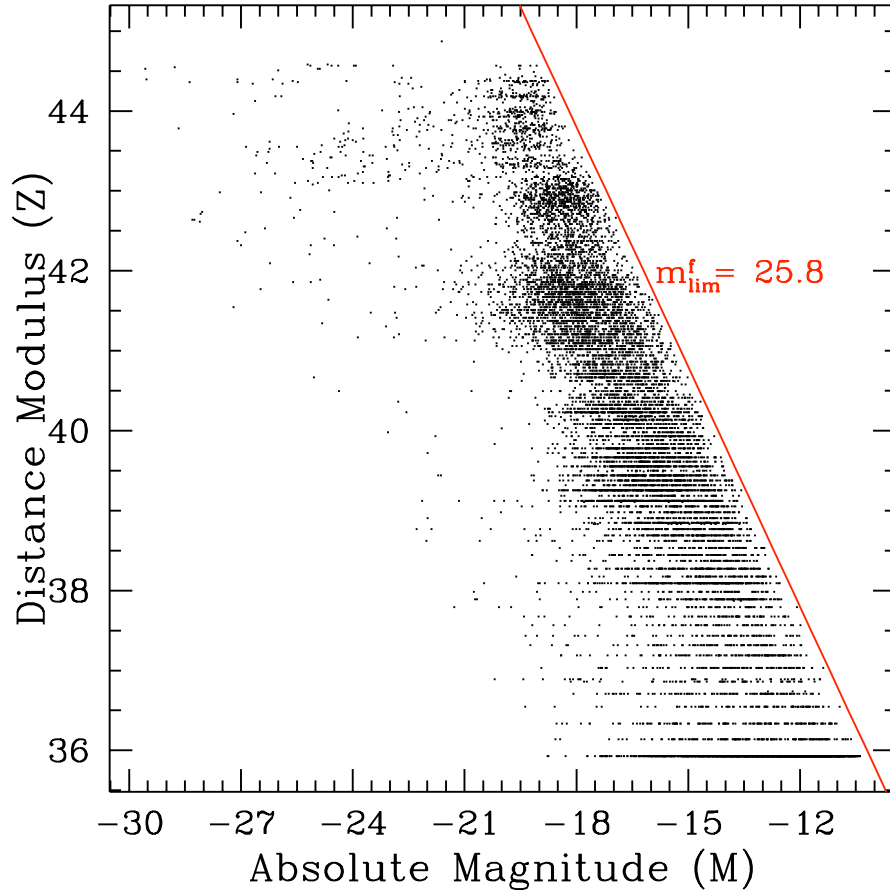


Figure 6.1: M - Z distribution of the CCLQGS data-set. The red line indicates the faint apparent magnitude cut-off of the sample, $m_{\text{lim}}^f = 25.8$ mag. Also evident is the low precision level of the photometric redshifts which manifests itself as discretisation on the M - Z plane.

therefore it was, at first, somewhat baffling. This was compounded since both statistics up to $m_* \approx 18.0$ mag had shown close agreement. It should also be noted that the case where T_c and T_v drop sharply as the limit of the survey has been passed, is due to the decreasing number of galaxies within the respective S_2 and S_4 regions. Similarly this implies that an observed sharp increase *above* 3σ would indicate an over-densed region.

The next step in our analysis, as in the case of the 2dFGRS, was to test any significant change in this result when adopting a bright limit. Figure 6.3 shows varying values of δZ and δM for T_c and T_v respectively in the range $0.1 < \delta Z, \delta M < 2.0$. What is immediately evident is the T_v spike remains for every value of δM and δZ . At this point it was becoming increasingly obvious (essentially by a process of elimination) that the precision level of the photometric redshifts coupled with the distribution of

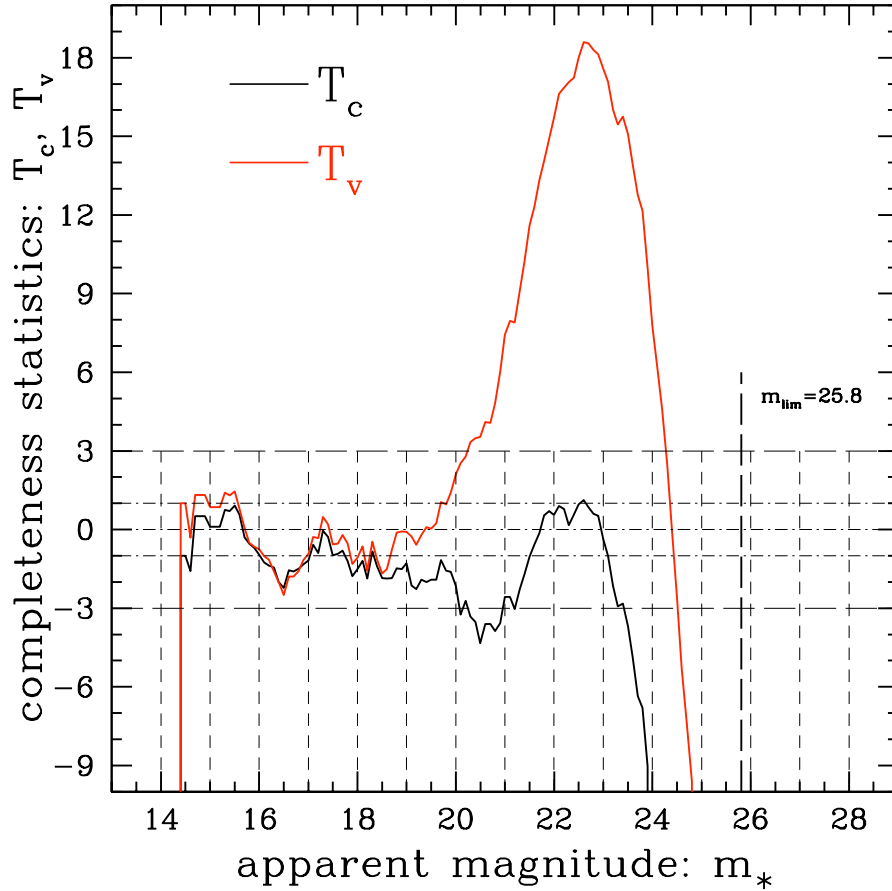


Figure 6.2: Initial T_c and T_v results of the CCLQGS data. We have assumed initially that there is no well defined bright limit and as such applied the original R01 statistic. Although both T_c and T_v fall below the -3σ before the official magnitude limit of $m_{\text{lim}} = 25.8$ mag, it is the distinct departure in completeness between the two statistics which peaks at around 18.5σ for T_v .

galaxies in the survey was playing a more crucial role than first thought.

Although never examined in detail, it was inferred in the [Rauzy \(2001\)](#) paper that finite precision data, in this case apparent magnitudes, will introduce discretisation to the data creating artificial gaps, and ‘steps’ in the magnitude distribution function. This in turn could introduce ‘spurious variations of the T_c statistic’. However, a well-documented way to overcome this effect requires the addition of a small amount of random noise or ‘jitter’ (e.g. see [Sivia and Skilling, 2006](#), p186) to the data to impose a ranking which therefore breaks any statistical ties. In the case of [Rauzy \(2001\)](#) a uniform random distribution between $[-0.005, 0.005]$ was added to the SSRS2 magnitudes to overcome any potential rounding problems.

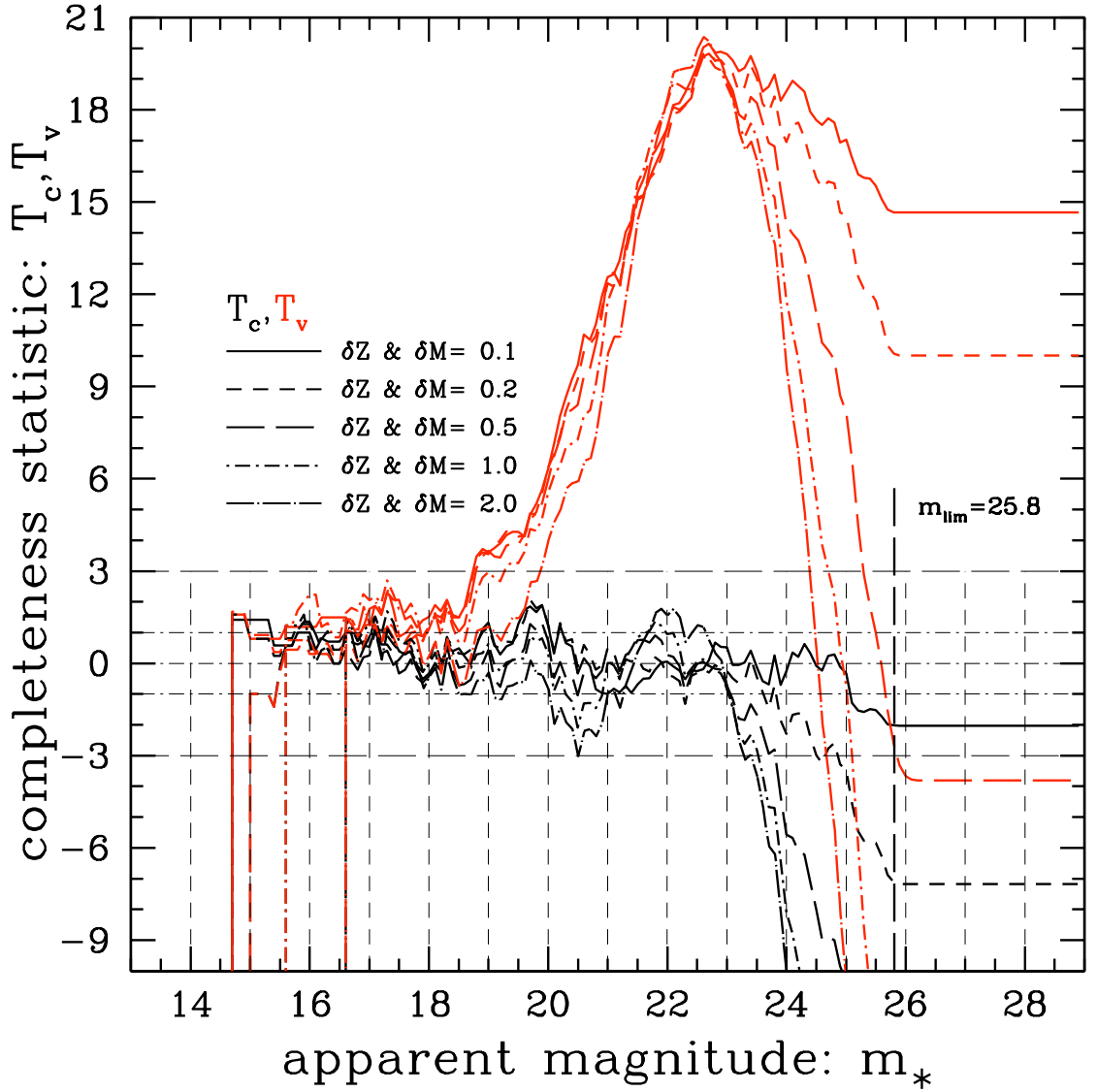


Figure 6.3: Initial T_c and T_v results of the CCLQGS data. We have assumed initially that there is no well defined bright limit and as such applied the original R01 statistic. Although both T_c and T_v fall below the -3σ before the official magnitude limit of $m_{\text{lim}} = 25.8$ mag, it is the distinct departure in completeness between the two statistics which peaks at around 18.5σ for T_v .

With this in mind, we examine more closely the M - Z distribution in Figure 6.1 and observe that the redshifts are rounded to 0.001 and appear heavily ‘quantised’ in the figure. Since our method considers the distance modulus, Z , which is essentially the $\log_{10}(z)$ of the redshift distribution, the quantised nature of the in Z is more apparent for nearby with lower values of z . We therefore adopt the same approach as in [Rauzy](#)

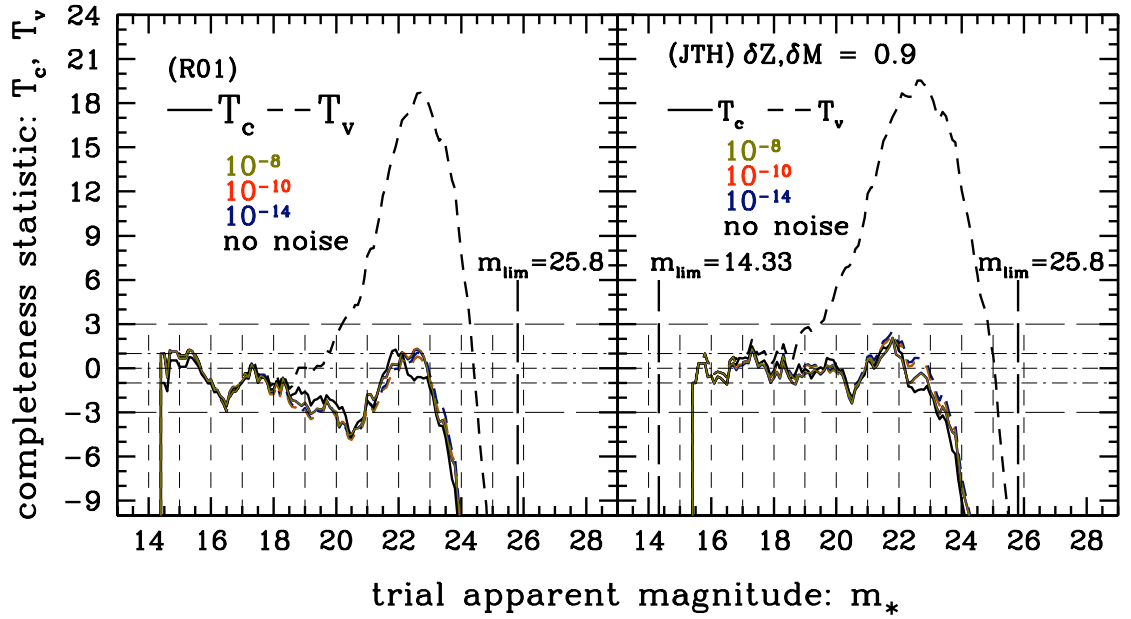


Figure 6.4: Results for the CCLQGS T_v statistic with random noise added to the raw redshift distribution. We add successive small amounts of trial uniformly distributed noise to the photometric redshifts between $[-10^{-8}, 10^{-8}]$, $[-10^{-10}, 10^{-10}]$ and $[-10^{-14}, 10^{-14}]$. The left-hand panel are results using the R01 method, whilst the right-hand panel shows the resulting JTH method. We observe that the T_v curve in all cases (red lines) shows no evidence of the spurious spike we originally found, and now shows similar behaviour as the T_c curve.

(2001) and apply different levels of uniformly distributed random noise to z denoted by an η parameter. In Figure 6.4 we have added an amount, η , of uniform random noise between $[-10^{-8}, 10^{-8}]$, $[-10^{-10}, 10^{-10}]$ and $[-10^{-14}, 10^{-14}]$ to z and compared it to the case of using the raw redshifts. The left hand panel in Figure 6.4 shows the R01 T_c and T_v statistics where the red solid line represents T_v with no noise added. The remaining red dashed curves represent $(z+\eta)$. As we can see, as soon as the statistical ties are broken by the addition of a small amount of random noise, the spike of T_v now reverts to the form of T_c . Similarly, in the right hand panel of Figure 6.4 we observe the same behaviour when we apply the Johnston, Teodoro and Hendry (2007) (hereafter referred to as, JTH) T_c and T_v statistics, adopting a bright magnitude limit of $m_{\text{lim}}^b = 14.33$ mag based on the brightest galaxy in the set.

Although we had established the cause of the spurious spike observed in T_v was a direct result of truncation of the redshifts, this did not completely explain why T_v departs from T_c at around $m_* = 18.5$ mag. The answer can be found by examining the M - Z distribution once again. In Figure 6.5 we have plotted this distribution and

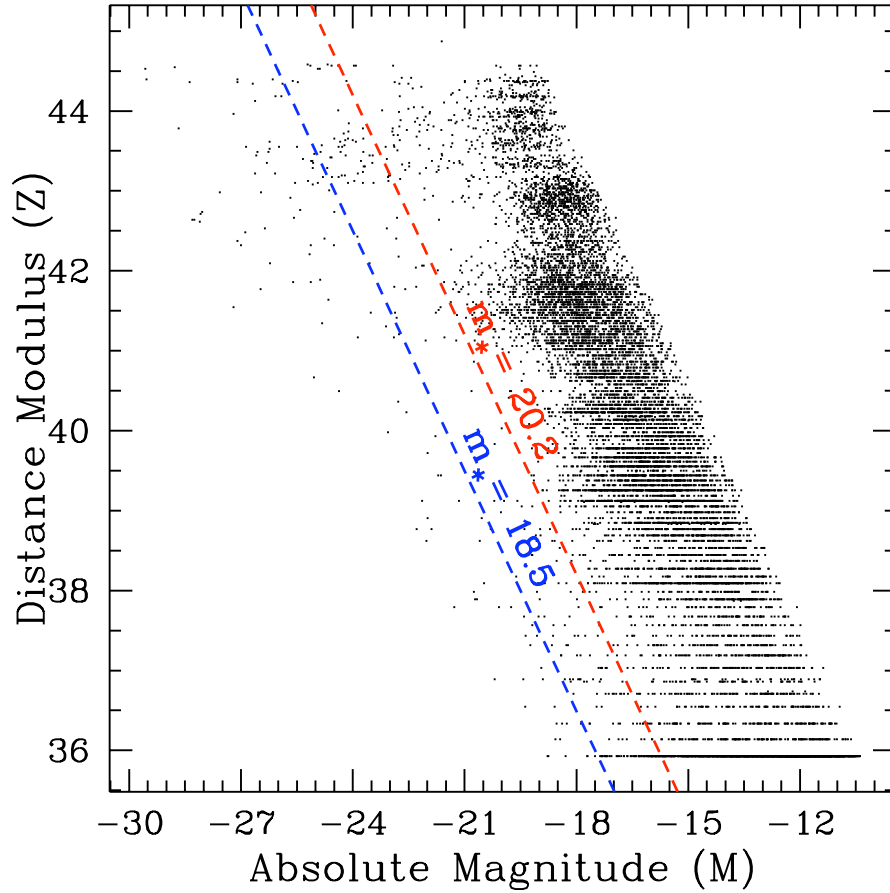


Figure 6.5: M - Z distribution for CCLQGS illustrating where T_v departs from T_c in Figure 6.3. The blue dashed line represents the trial limit magnitude, $m_* = 18.5$ mag and is approximately where we see T_v initially departing from T_c . The red dashed line is for $m_* = 20.2$ mag and is approximately where T_v first crosses the 3σ confidence limit on its way to peak at $m_* \approx 22.5$ mag.

superimposed two key trial magnitude limits, m_* - the first, where T_v initially departs from T_c (the blue dashed line), and the second where T_v crosses the 3σ confidence limit at $m_* \approx 20.2$ (the red dashed line). Up to $m_* = 18.5$ mag, the distribution is sparsely populated by distant bright galaxies (top left of the M - Z diagram in Figure 6.1) which show little in the way of discretisation due to the fact we are observing these galaxies in terms of their distance modulus. As we move beyond $m_* = 18.5$ mag to $m_* = 20.2$ mag we can see the nearby galaxies are now playing a more crucial role in the T_v calculation where the discretisation in Z is now clearly visible. This suggests that the more ‘quantised’ galaxies at approximately $Z < 40$ dominate this truncation effect observed in T_v .

To test this we used the raw redshifts and simply sampled slices of the data and applied our statistics in two opposite ways using the JTH method and adopting a δZ and $\delta M = 1.0$ throughout. The first approach applied T_c and T_v in Figure 6.6-top, starting from the furthest galaxies in the survey and working in decreasing limiting values of redshift (as shown in the top right-hand panel). Thus, the ranges in redshift we considered were $1.0 < z < 2.5$, $0.8 < z < 2.5$, $0.4 < z < 2.5$ and $0.1 < z < 2.5$. The top left-hand plot of Figure 6.6 shows the resultant T_c and T_v curves for each range in redshift. What is immediately clear is that between $0.4 < z < 2.5$, T_v shows none of the spiking characteristics seen when testing the whole data-set. Even when we extend the computation back to a $z_{min} = 0.1$ there is only a relative marginal peak present in T_v at $\sim 6.5\sigma$. We therefore, test the data in the opposite direction by starting at the minimum redshift, $z_{min} = 0.05$, and sampling at increasing maximum limits of redshift: $0.05 < z < 0.15$, $0.05 < z < 0.2$, $0.05 < z < 0.5$ and $0.05 < z < 1.0$ as illustrated on the bottom right-hand panel of Figure 6.6. The solid red line on the bottom left-hand panel of Figure 6.6 represents the T_v results for the $0.05 < z < 0.15$ slice and shows a strong peak at $T_v \sim 33\sigma$. As we move to increasing limiting values of redshift we observe this peak fall towards the level shown in the right-hand panel of Figure 6.4. Therefore, our results are indicating that the relatively nearby galaxies $z \lesssim 0.15$ are dominating the atypical nature of this bias.

The reason why we observe this effect only in T_v and not in T_c is now straightforward to explain. Figure 6.7 shows the construction of both T_v (left) and T_c (right) for a galaxy at (M_i, Z_i) . For the construction of the regions S_3 and S_4 for T_v we can see that for all galaxies lying at the same distance modulus as Z_i will always be counted in S_3 by the definition given in Equation 5.4 on page 97. Since there are fewer nearby galaxies in distance modulus this introduces a bias into the estimator which can only be overcome by either adding random noise to the redshift data to break the ties, or improving on photometric redshift model predictions to begin with (for example see Sheth, 2007). It becomes apparent that for the case of rounded data in redshift, the T_c will not be affected, since the construction of this statistic samples the cumulative distribution of absolute magnitudes and not redshifts. Therefore, it follows that if the magnitudes suffered from a similar truncation one would expect to see an adverse effect in the T_c statistic but not necessarily with T_v .

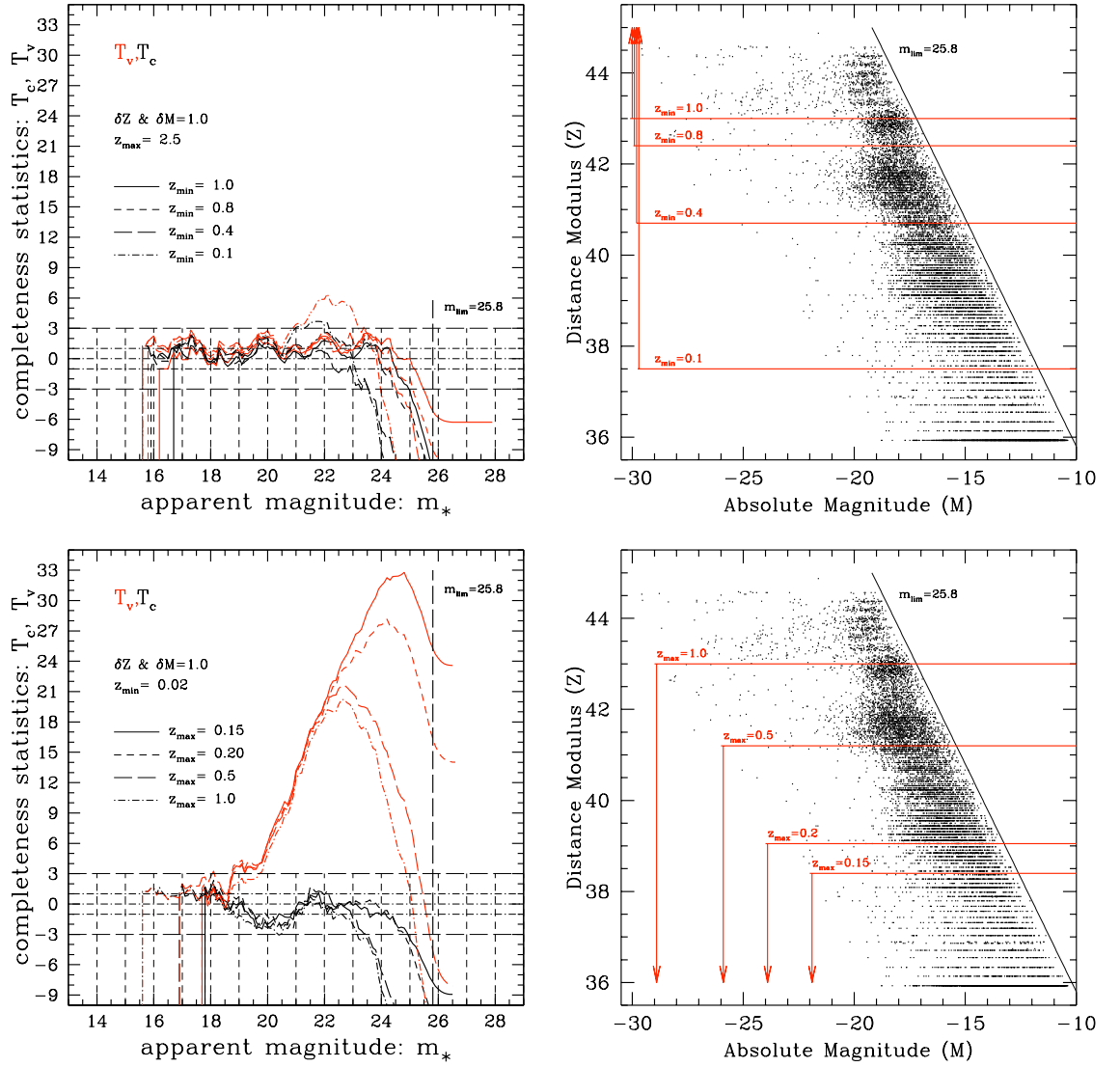


Figure 6.6: The top-left panel represents the JTH T_c and T_v curves for four slices in redshift, z . We have adopted a δZ and $\delta M = 1.0$ throughout. The top-right hand panel shows schematically how the redshift slices: $1.0, 0.8, 0.4$ and $0.1 < z < 2.5$, relate in the context in the M - Z plane. The completeness statistics clearly show that within the range of z_{\min} slices from $z_{\min} > 0.4$, both T_c and T_v behave as one would expect for a complete sample. As we move to a $z_{\min} = 0.1$ we see the strange feature in T_v begins to reappear. The bottom panel on the left represents the resulting T_c and T_v curves for the following four slices in redshift: $0.05 < z < 0.15$, 0.05 , 0.5 and 1.0 , and is shown schematically in the bottom right-hand panel. Both completeness statistics in this case clearly show a much greater peak that we have observed previously within the small range $0.15 > z_{\min} > 0.02$. As we move to greater limiting redshifts we observe the peak reverting to its expected confidence level for a complete survey.

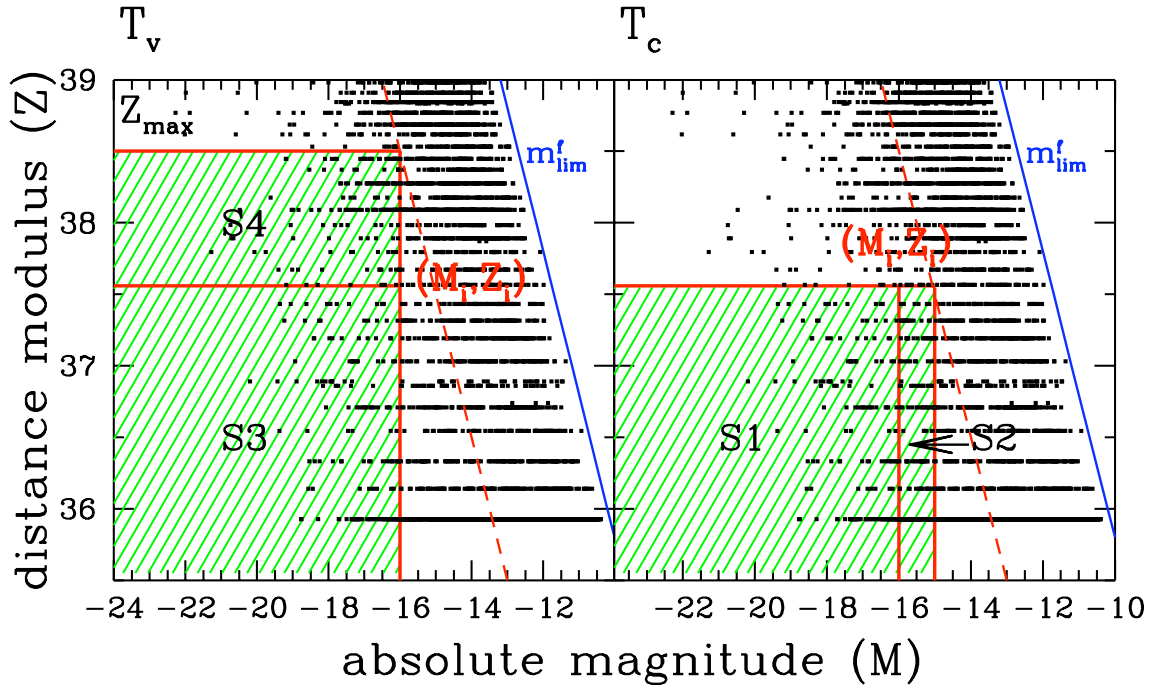


Figure 6.7: Construction of T_c and T_v to illustrate the bias introduced into the T_v estimator due to truncated redshifts. The left and right-hand panels show the familiar construction for T_v and T_c respectively, where we do not include the presence of a bright limit.

6.3 Final Completeness Assessment of the Data

Having established the causes surrounding the large spike observed in T_v , we continue with our completeness analysis. For the remainder of this chapter we will be adding an amount, η , of uniformly distributed random noise between $[-10^{-8}, 10^{-8}]$ to the raw redshift data.

Let us re-visit Figure 6.2 on page 103 where the R01 T_c method has been applied. The curve firstly shows a dip in completeness with a $T_c \sim -4.7\sigma$ minimum at a limiting apparent magnitude $m_* \sim 20.5$. There is then a slight recovery within $|T_c| < 3\sigma$ before dropping systematically and indicating a *true* apparent magnitude limit of $m_* \sim 23.2$ mag, an apparent magnitude which is considerably brighter than the limit of the survey data of $m_{\text{lim}} = 25.8$ mag. This behaviour is not unlike the initial analysis of the 2dFGRS in Figure 4.3 on page 69. By adding the noise element to the redshifts, Figure 6.4-left showed that T_v follows a similar trend to T_c . In the case of 2dFGRS, we concluded that a bright limiting magnitude would have to be included into our completeness calculation in order to make an accurate assessment of the data. Whilst examining the T_v ‘spike’ in the previous section we adopted a bright limit which, as

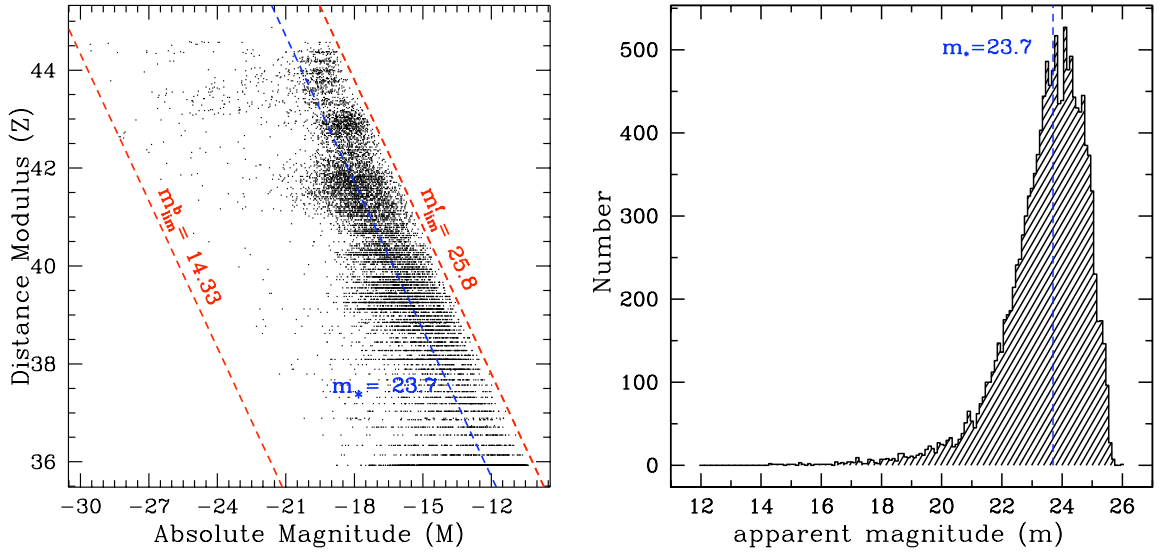


Figure 6.8: M - Z distribution for the CCLQGS data showing the average completeness magnitude limit derived from T_c and T_v . The blue dashed line on the left-hand plot shows the magnitude limit as identified by both statistics, $m_* = 23.7$ mag. The right-hand plot is a histogram of the apparent magnitudes for the same sample with the same limiting magnitude superimposed shown as the blue dashed line. There is a visible turnover in the magnitudes between $23 \lesssim m \lesssim 25.0$ indicating a steady drop in the number of galaxies towards the magnitude limit instead of a sharp cut in magnitude at the limit. This, however, does not explain why both statistics have systematically dropped below the -3σ limit by $m_* = 23.7$ mag.

with the 2dFGRS, was chosen as the brightest galaxy in the sample $m_{\text{lim}}^b = 14.33$ mag. We then applied the JTH statistic for a selection of δZ and δM between $[0.1, 2.0]$ (see Figure 6.3 and Figure 6.4-right). The results show that the initial dip below -3σ is no longer there, however, as we increase the size of both δZ and δM , the completeness results for T_c and T_v indicate a varying range for the *true* apparent magnitude limit an effect that we examine in more detail in the following chapter. Moreover, if we look once again at Figure 6.6 we also observe a change in the magnitude limit for increasing slices in redshift.

What now remains is the question - why do T_c and T_v indicate an average minimum magnitude limit of $m_{\text{lim}} \sim 23.7$ mag which lies at a point considerably brighter than the limit of the survey data? Unfortunately, at this time there is no immediate answer, and for the moment we can only speculate. In the right-hand panel of Figure 6.8 we have plotted the apparent magnitude distribution of our sample with the superimposed blue line indicating our completeness magnitude limit result. It is obvious in this plot that the distribution is not defined by a sharp cut-off, but instead turns over gradually at

approximately where our method is indicating the actual limit of the survey. However, this is inconsistent with what we would expect in such a scenario. If a data-set is defined with an m_{lim} that has a gradual tail off, then as a result we would expect to see T_c and T_v drop gradually to reflect this. However, our statistics drop systematically *before* the turnover in the magnitude distribution begins.

6.4 Conclusions

The completeness analysis of this survey sample has prompted us to consider possible limitations of the method and/or the survey catalogue that may not have surfaced otherwise. When applying the T_v estimator we observed a very distinct and large peak which had not been observed in the analysis of previous surveys. We concluded that this effect was a direct result of highly rounded photometric redshift data that, in turn, created large artificial gaps, and ‘steps’ in the distribution function for the distance modulus for the relative nearby galaxies. Consequently, this introduced a strong bias in the T_v estimator since the range in redshift for which this effect dominates accounts for approximately one third of the total number of galaxies in data-set. This led us to add a small amount of statistical ‘jitter’ to the redshift data to break any statistical ties that may be adversely affecting our calculation. This additional step appears to have corrected this issue.

Although, this procedure is a well recognised one within the statistics community, it would be too easy to be judgmental of the T_c and T_v estimators (and indeed other estimators that this may affect) as it brings to light a key issue regarding the precision of the photometric redshift measurements. The addition of any amount of random noise to data is essentially adding physical precision to an instrument, and therefore, the model predictions, that have clear limitations in measurements.

Our completeness results also indicated a possible range for the true faint magnitude limit as we varied the values of δZ and δM in the JTL method. This, combined with a similar result from the 2dFGRS, has lead us to propose a more efficient and optimised way to compute both T_c and T_v that should provided the user with a confidence level for the choice of m_{lim} based on the signal-to-noise of the system. We explore this in the following chapter.

Finally it is important to reiterate that the magnitudes presented here have not yet been corrected for any form of evolution, that will undoubtedly be inherent with such a deep redshift survey. Nor have they been k -corrected. These are both crucial factors which would potentially have strong impact on the completeness assessment for

a survey that extends to such high redshifts. Prior to the testing of this sample we have only examined data-sets which are relatively shallow in redshift, where the effects of evolution are less dominant. Since we observe a change in the magnitude limit assessment for different slices in redshift may indeed suggest evolutionary effects are impacting on our statistics. Future studies could include simulating such surveys out to high redshifts such that we can control and explore the impact of all these effects.

Chapter 7

Optimising T_c and T_v

“Science is always wrong; it never solves a problem without creating ten more.”

George Bernard Shaw

In this chapter we present our ongoing research into an area concerning the optimisation of the generalised JTH T_c and T_v statistics. The need for this research has been prompted largely by the completeness results from the 2dFGRS in chapter 4 and the CCLQG data in the previous chapter. In both cases we briefly discussed that by varying the widths of δZ and δM two distinct side-effects for the determination of the true m_{lim}^f were revealed:

1. For very small values of δZ and δM the respective T_c and T_v statistics will be dominated by shot-noise, making it impossible to draw statistical conclusions concerning the true faint apparent magnitude limit, m_{lim} .
2. Conversely, for values of δZ and δM that are very large we observe a range in possible values for the faint magnitude limit for data-sets that are not well described by a sharp m_{lim} .

Therefore, we now examine the 2dFGRS and CCLQG data-sets in more detail along with the other surveys we previously examined, namely MGC and SDSS (Early Types). We demonstrate that by choosing to estimate the random variables with fixed values of δZ and δM we unavoidably introduce the above effects which, if not accounted for properly, could adversely influence our conclusions with regards to our estimation of

the correct faint magnitude limit, m_{lim}^f . In final part of this chapter we will outline an alternative approach to estimating the random variables that amounts to a direct determination of the signal-to-noise of the system and which should circumvent issues which we now describe.

7.1 Current Issues with the ζ and τ Estimators

Firstly, let us recall our definitions for the random variables ζ and τ for respective statistics T_c and T_v ,

$$\zeta = \frac{F(M) - F[M_{\text{lim}}^b(Z - \delta Z)]}{F[M_{\text{lim}}^f(Z)] - F[M_{\text{lim}}^b(Z - \delta Z)]} = \frac{S_1}{S_1 \cup S_2} = \frac{r_i}{n_i + 1}, \quad (7.1)$$

$$\tau = \frac{H(Z) - H[Z_{\min}(M - \delta M)]}{H[Z_{\max}(M)] - H[Z_{\min}(M - \delta M)]} = \frac{S_3}{S_3 \cup S_4} = \frac{q_i}{t_i + 1}. \quad (7.2)$$

where the points for each region are represented by, r_i belonging to S_1 , n_i belonging to $S_1 \cup S_2$, q_i belonging to S_3 , and t_i belonging to $S_3 \cup S_4$. Essentially, the milestone of our extension to the R01 method lay in the introduction of the fixed quantities δZ for ζ and δM for τ . Fixing these quantities to a predetermined width allowed us to re-construct the regions in Equations 7.1 and 7.2 within any doubly truncated survey i.e. for a survey with well defined bright and faint apparent magnitude limits. Let us firstly examine the results obtained from all the of the above surveys for the case of small values of δZ and δM .

7.1.1 Shot noise - the ‘flat-line’ effect

In the top two panels of Figure 7.1 we have applied T_c and T_v to the 2dFGRS over the range $0.001 < \delta Z, \delta M < 0.008$ in increments of 0.001. We observed that all curves for both test statistics fluctuate within the 3σ limits for each trial m_* within $13.346 < m_* < 19.45$. However, once m_* moves beyond the published limit of the survey all the curves drop slightly and then flatten or ‘flat-line’ inside $-3\sigma < T_c, T_v < 3\sigma$ instead of dropping sharply below the -3σ level as one, by now, might expect. Similarly, the top panels of the Figures 7.2, 7.3 and 7.4 demonstrate the same flat-lining effect for the SDSS (Early Types), MGC and CCLQG respectively. For the SDSS-Early Types, observe that for values $0.001 < \delta Z, \delta M \lesssim 0.01$, T_c and T_v flat-line beyond the survey limit of $m_{\text{lim}}^f = 17.55$ within the 3σ limits. As we move to increasing values of δZ and δM as shown in the bottom panels, the curves remain flat-lined beyond the magnitude

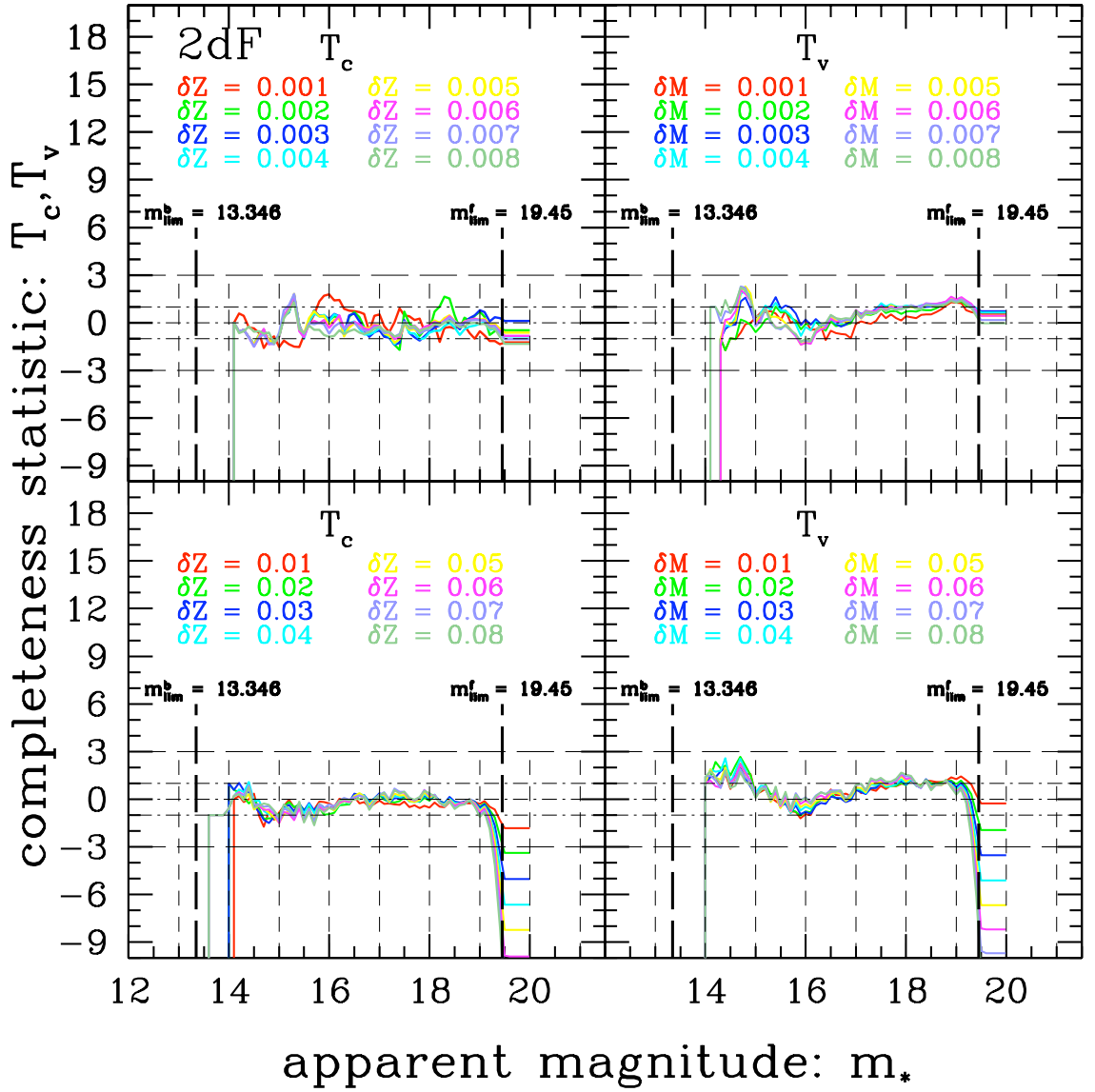


Figure 7.1: T_c and T_v results for the 2dFGRS applying the Johnston et al. (2007) generalisation for values in the range of $0.001 < \delta Z, \delta M < 0.08$. We observe for values $0.001 < \delta Z \lesssim 0.01$ and $0.001 < \delta M \lesssim 0.02$ T_c flat-lines within the 3σ limits thus indicating that the noise level is greater than the ζ and τ signal. However, as we move to increasing values of δZ and δM as shown in the bottom panels, the curves remain flat-lined beyond the magnitude limit but are now below -3σ indicating that there is enough signal in ζ and τ to draw meaningful statistical conclusions.

limit but are now below -3σ . For the MGC in the range $0.001 < \delta Z, \delta M \lesssim 0.06$, T_c and T_v flat-line within the $|3\sigma|$ limits. Once again, as we move to increasing values of δZ and δM as shown in the bottom panels, the curves remain flat-lined beyond the magnitude limit but are now below -3σ . Finally, the CCLQG survey in Figure 7.4,

confirms this behaviour as we observe both statistics flat-line between the same values as MGC, $0.001 < \delta Z, \delta M \lesssim 0.06$.

This flat-lining is essentially an indicator of the shot-noise level for the chosen width of δZ and δM and understanding why this happens is quite straightforward when we consider the contributing factors: the number of objects in the catalogue and the range in apparent magnitude, m , of the survey. Figure 7.5 illustrates how the shot-noise dominates for small values of δZ and δM . In our example we look at the 2dFGRS data-set. The top left panel shows the now familiar M - Z distribution with the red diagonal lines representing the faint apparent magnitude limit m_{lim}^f and our adopted bright limit m_{lim}^b . The main feature of this plot is the narrow red, blue and green lines which actually define the T_c regions S_1 and S_2 for a galaxy at (M_i, Z_i) for $\delta Z = 0.001$, 0.008 and 0.02 respectively, where we are considering a trial m_* equal to the survey limit i.e. $m_{\text{lim}}^f = 19.45$. Since these ‘strips’ represent such a tiny fraction of the diagram, the top right hand plot shows a close up of this particular region where the areas are clearly defined. The bottom panels in Figure 7.5 represents, for the same galaxy at (M_i, Z_i) , the T_v construction for $\delta M = 0.001$, 0.008 and 0.02 . What is immediately obvious for both T_c and T_v is the relative small number of galaxies that are counted within the S_1 and S_2 regions for each δZ and δM . By considering Equations 7.1 and 7.2, it becomes more obvious that as we move m_* beyond the m_{lim}^f cut-off, the T_c and T_v statistics will see the same small fractional change in the relative numbers between r_i and n_i and q_i and t_i respectively for each $m_* > m_{\text{lim}}^f$, and hence the shot-noise is dominant. In other words, the fact the flat-line occurs within the 3σ limit is perhaps an indication that the shot noise level is greater than the ζ and τ ‘signal’. Therefore, it also follows that as we increase δZ and δM in size, the number of galaxies counted within the stated regions also increase and thus the respective statistics become more sensitive to m_{lim}^f and the signal-to-noise also increases.

Finally, since we are dealing with flux-limited catalogues with a finite number of galaxies, we can show via Figure 7.6 that by applying the original Rauzy completeness test, which consequently has no restriction in the height of the constructed regions, and allowing m_* to move far enough beyond m_{lim}^f then ultimately we still observe shot-noise, albeit at a much larger negative level of σ .

7.1.2 Variation in m_{lim} effect

In this section we explore the apparent variation in determining the true m_{lim} resulting from larger values in δZ and δM of a survey that is doubly truncated. As we have

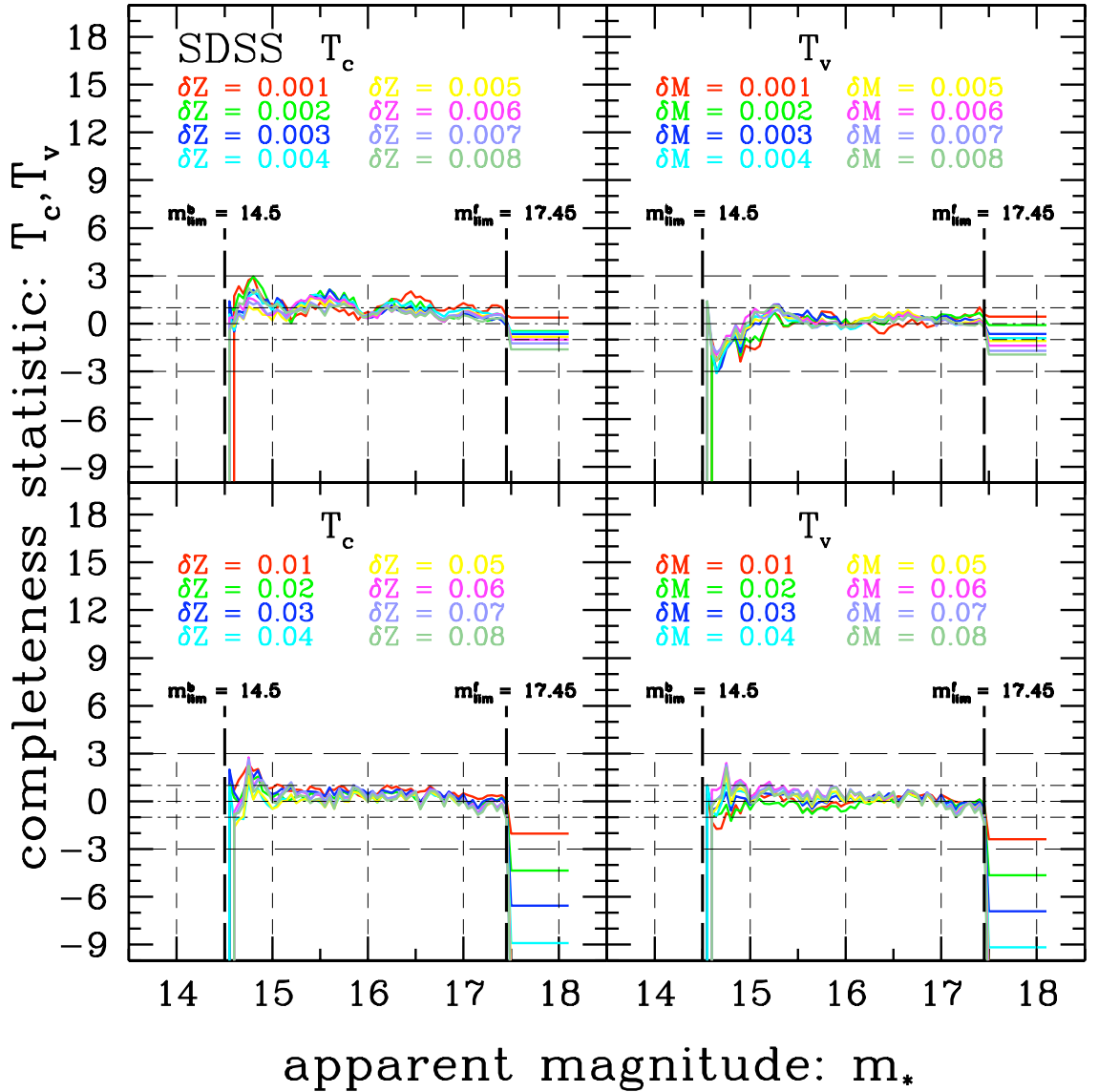


Figure 7.2: T_c and T_v results for SDSS applying the JTH generalisation for values in the range of $0.001 < \delta Z, \delta M < 0.08$. We similarly observe, as with the 2dFGRS, that for values $0.001 < \delta Z, \delta M \lesssim 0.01$, T_c and T_v flat-line within the 3σ limits indicating that the noise level is greater than the ζ and τ signal. Once again, as we move to increasing values of δZ and δM as shown in the bottom panels, the curves remain flat-lined beyond the magnitude limit but are now below -3σ .

seen from the previous section, if δZ and δM are sufficiently small then the resulting shot-noise associated with the T_c and T_v curves will dominate and therefore flat-line within the 3σ limits. This implies that one has to choose larger values such that both statistics will drop below -3σ once the *true* magnitude limit has been identified. However, as we shall demonstrate, this in turn has potential problems. If we now look

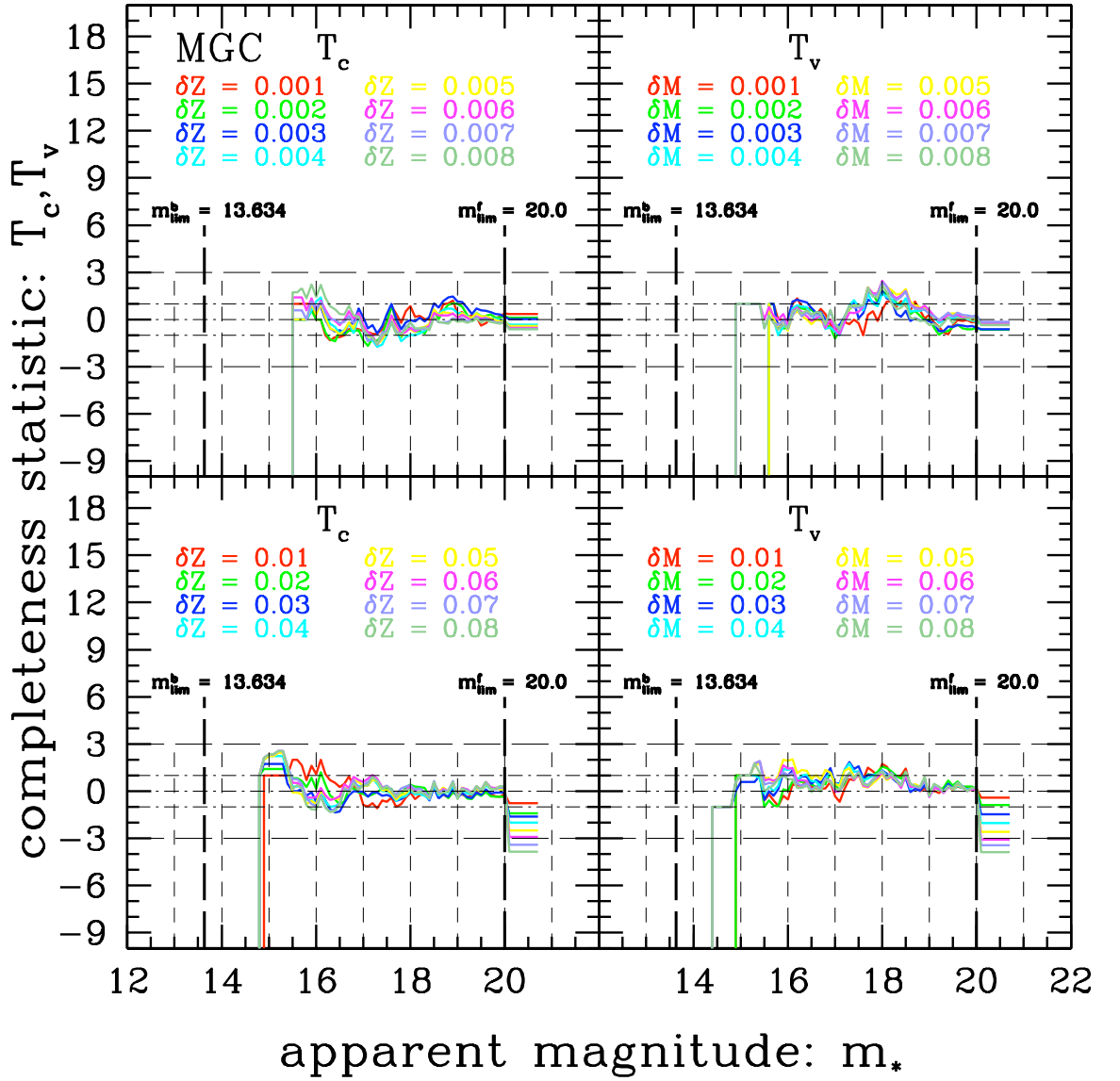


Figure 7.3: T_c and T_v results for MGC applying the JTH generalisation for values in the range of $0.001 < \delta Z, \delta M < 0.08$. We find that for values $0.001 < \delta Z, \delta M \lesssim 0.06$, T_c and T_v are flat-line within the 3σ limits indicating that the noise level is greater than the ζ and τ signal. Once again, as we move to increasing values for $\delta Z, \delta M > 0.06$ as shown in the bottom panels, the curves remain flat-lined beyond the magnitude limit but are now slightly below -3σ .

at Figures 7.7-2dFGRS, 7.8-SDSS and 7.9-CCLQG we can quite clearly see that as δZ and δM increase, the point at which the test statistics systematically fall below our -3σ limit (indicating the true apparent magnitude limit) varies with it. In the case of the 2dFGRS, on the interval $0.003 < \delta Z, \delta M < 3.0$ we actually observe a respective range of m_{lim} from $19.4 \lesssim m_{\text{lim}} \lesssim 19.0$. For SDSS on the interval $0.4 < \delta Z, \delta M < 2.0$

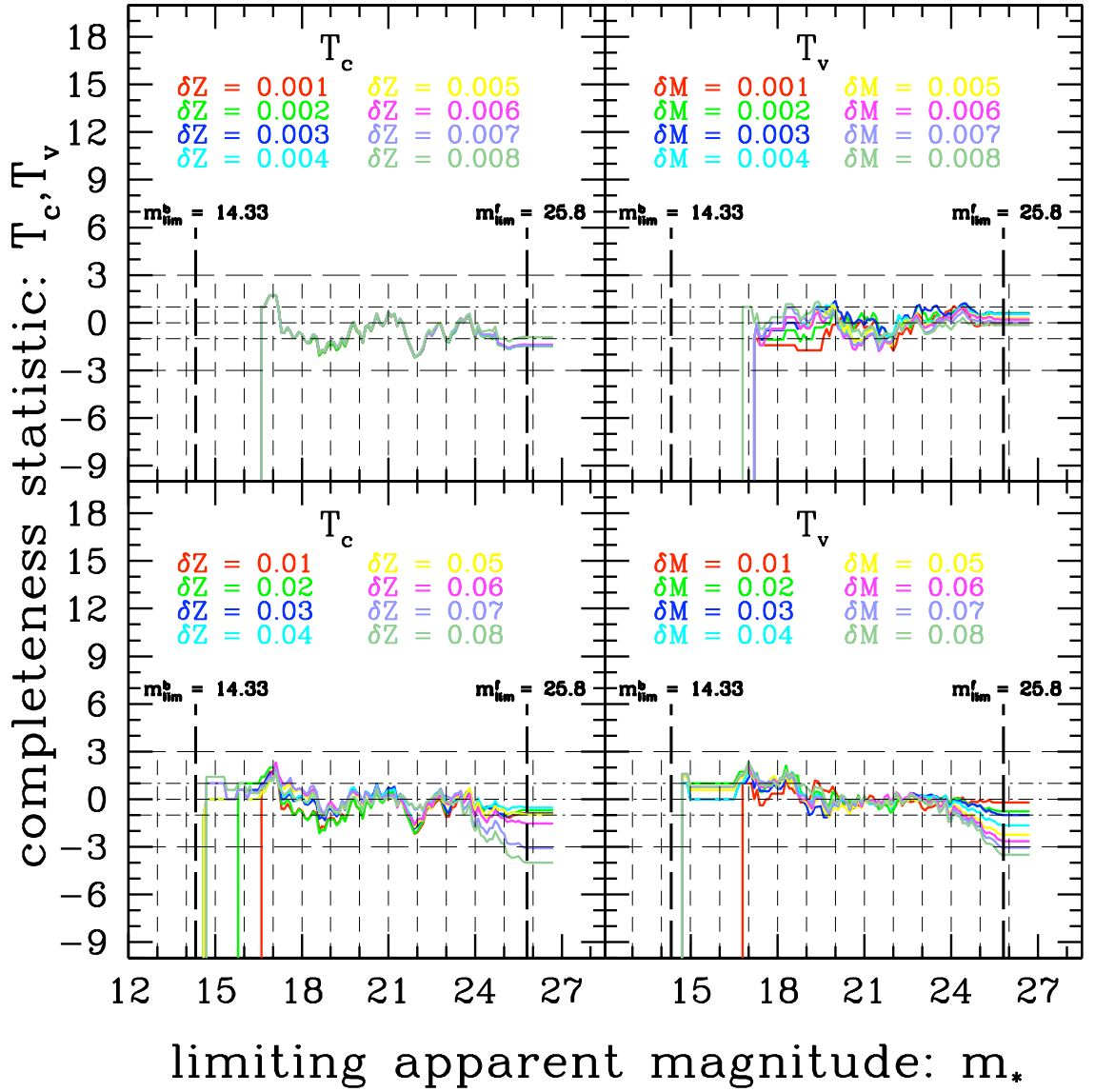


Figure 7.4: T_c and T_v results for the CCLQG data between $0.001 < \delta M, \delta Z < 0.08$. The top panels show T_c and T_v in the range $0.001 < \delta M, \delta Z < 0.008$ where we observe the flat-lining effect. As we move to $0.01 < \delta M, \delta Z < 0.08$, we observe the same behaviour as with the 2dFGRS and SDSS-Early Types. In this case the T_c curves begin to drop below -3σ for values of δM and $\delta Z = 0.07$ and 0.08 corresponding to $m_* = 25.8$ mag and $m_* = 25.6$ mag respectively. The T_v curves, however, show for δM and $\delta Z = 0.07$ and 0.08 , limiting magnitudes of $m_* = 25.8$ mag and $m_* = 25.1$ mag respectively.

the respective observed shift is $17.15 \lesssim m_{\text{lim}} \lesssim 17.45$.

For the CCLQG survey (Figure 7.4) we see both T_c and T_v flat-line within the 3σ limits for $0.001 < \delta M, \delta Z < 0.06$ once m_* has moved beyond the limit of the survey. We then find for the range $0.07 < \delta M, \delta Z < 0.8$ both estimators indicated completeness for

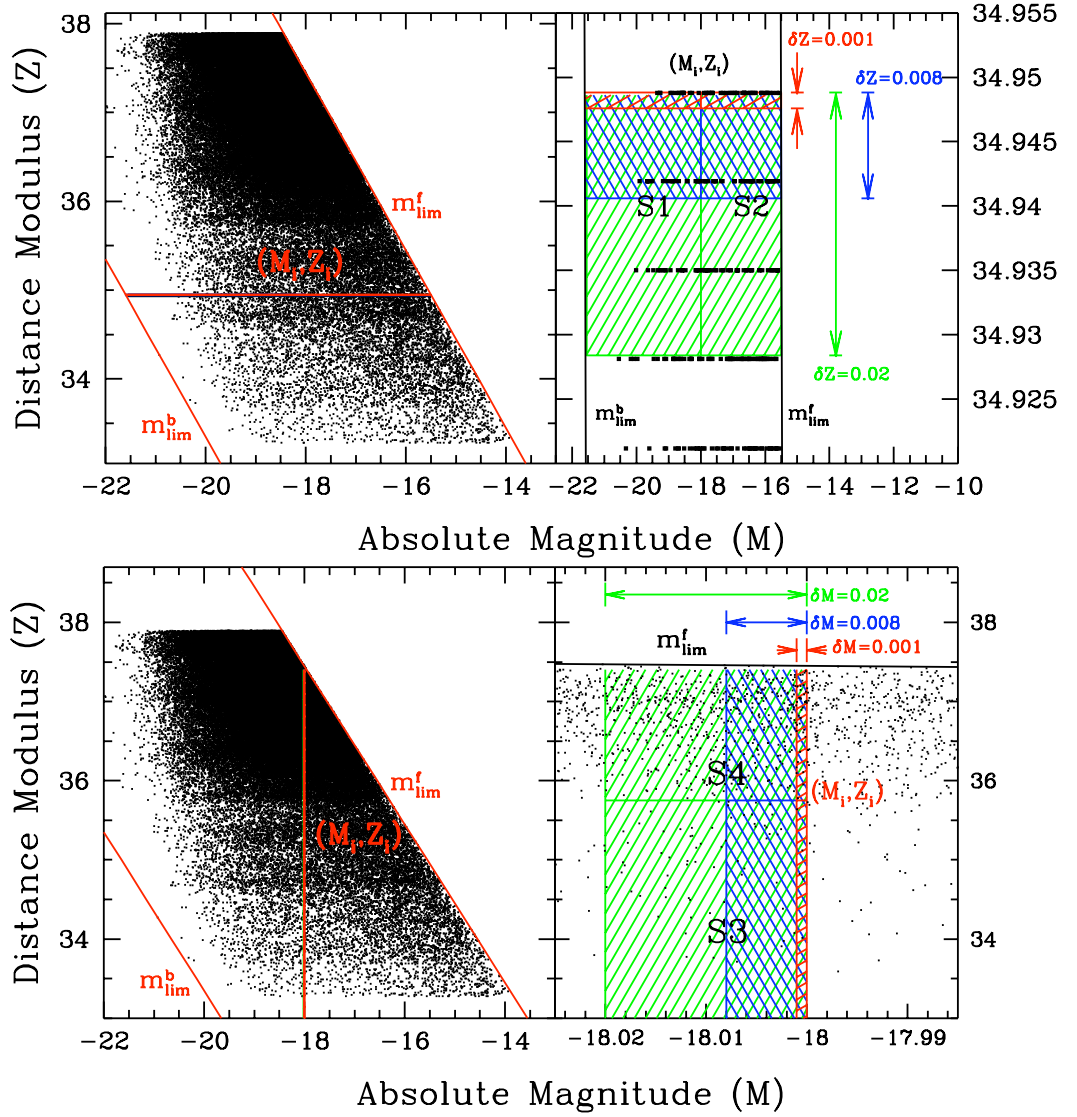


Figure 7.5: Schematic illustrating the cause of the flat-lining effect (within the 3σ confidence limits) observed for small values of δZ and δM . The left hand panels, both top and bottom, show the (M, Z) distribution for the 2dF with the faint apparent magnitude limit m_{lim}^f and our adopted bright limit, m_{lim}^b shown as red diagonal lines. The top-left plot considers the ζ construction for a galaxy at (M_i, Z_i) , with $\delta Z = 0.001, 0.008$ and 0.02 . The top-right plot is a zoomed in version of the left to allow us to see the three distinct regions created by the increasing sizes of δZ and the relative number of galaxies contained therein. Similarly, the bottom left panel shows the τ construction for a galaxy at (M_i, Z_i) for δM values that are equal to δZ . Consequently, the bottom-right panel is the zoomed in version of the bottom left.

the survey between the respective range of limiting magnitudes, $23.7 < T_c(m_*) < 25.6$ and $23.7 < T_v(m_*) < 25.8$. The bottom panels on Figure 7.9 show that $1.0 < \delta M, \delta Z < 4.0$ then converge to a $T_c(m_{\text{lim}}^*), T_v(m_{\text{lim}}^*) \sim 23.6$ mag, a behaviour also echoed in the

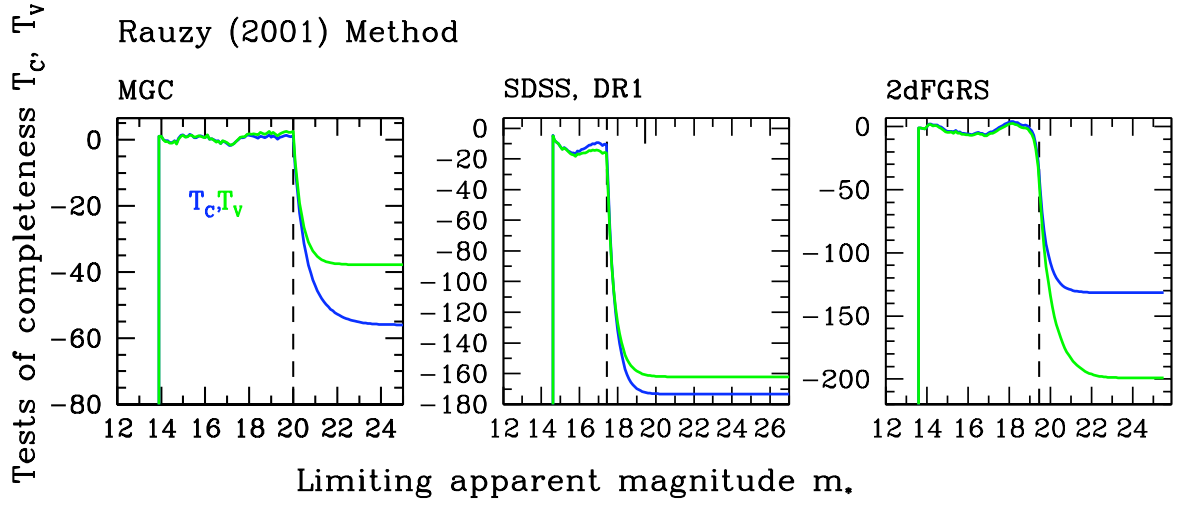


Figure 7.6: T_c and T_v plots for MGC, SDSS and 2dFGRS using the R01 method to illustrate the flat-line effect. We can see in all three plots that if one allows m_* to pass far enough beyond the magnitude the shot-noise level will eventually dominate albeit for extremely negative values of σ .

2dFGRS analysis.

However, as one might expect, the MGC completeness results shown in Figure 7.10 for $0.1 < \delta Z, \delta M < 4.0$ show *no* variation in m_{lim} and show a consistent drop in both statistics at $m_{\text{lim}} = 20.0$. This result is not unexpected since the application of the R01 method to this survey in § 3.2 yielded the same result.

This scope for variation in identifying the true magnitude limit of a survey could potentially lead anyone applying these methods to select the ‘best’ result possible for their data-set without knowing or being aware of the actual uncertainty of their result. This would then in turn defeat the purpose of the original Rauzy completeness test where one should be able to identify non-parametrically, and therefore validate independently, the completeness level of a given data-set up to a given magnitude limit.

We therefore must consider at this stage the possibility that neither 2dFGRS, SDSS or CCLQG are complete up to their respective apparent magnitude limits of $m_{\text{lim}}^f = 19.45$, $m_{\text{lim}}^f = 17.55$ mag, and $m_{\text{lim}}^f = 25.8$ and have to provide an alternative route to estimating our random variables which allow the user to estimate the error for each T_c and/or T_v point and simultaneously choose a signal-to-noise threshold that allows a varying δZ and δM to determine the true limit of the survey independent of the potential bias from the user.

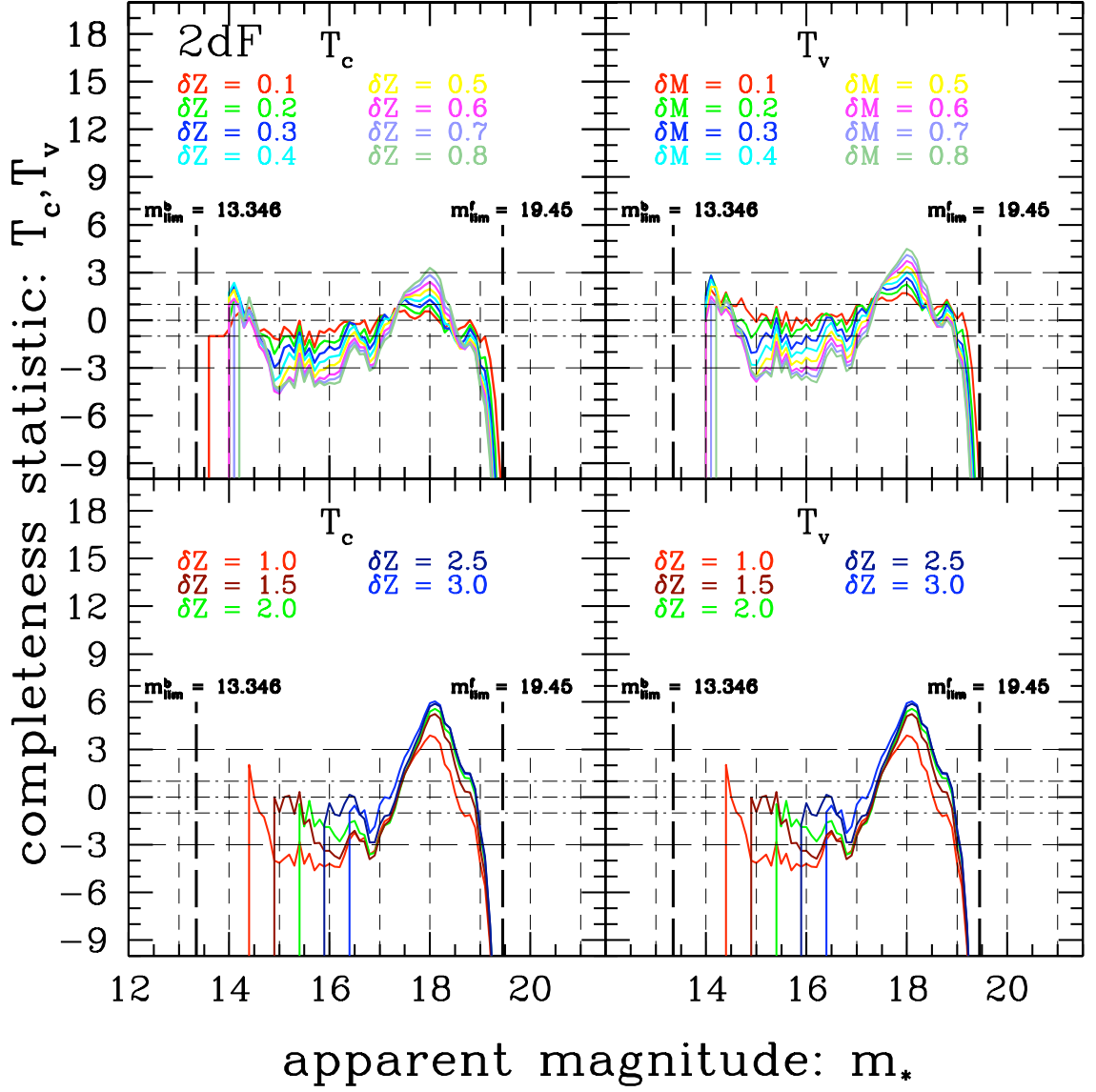


Figure 7.7: T_c and T_v results for 2dFGRS applying the JTH method for values in the range of $0.1 < \delta Z, \delta M < 3.0$. In all four panels we observe that as both δZ and δM increase, the T_c and T_v curves resemble the characteristic shape when we applied the R01 method which assumes a faint apparent magnitude limit only. Moreover, the incremental increase in the random variables results in systematic shift in the magnitude limit where T_c and T_v cross the -3σ limit. If we include results from Figure 7.1 such that $0.003 < \delta Z, \delta M < 3.0$, we observe then that our resulting trial m_{lim} ranges from $19.4 \lesssim m_{\text{lim}} \lesssim 19.0$.

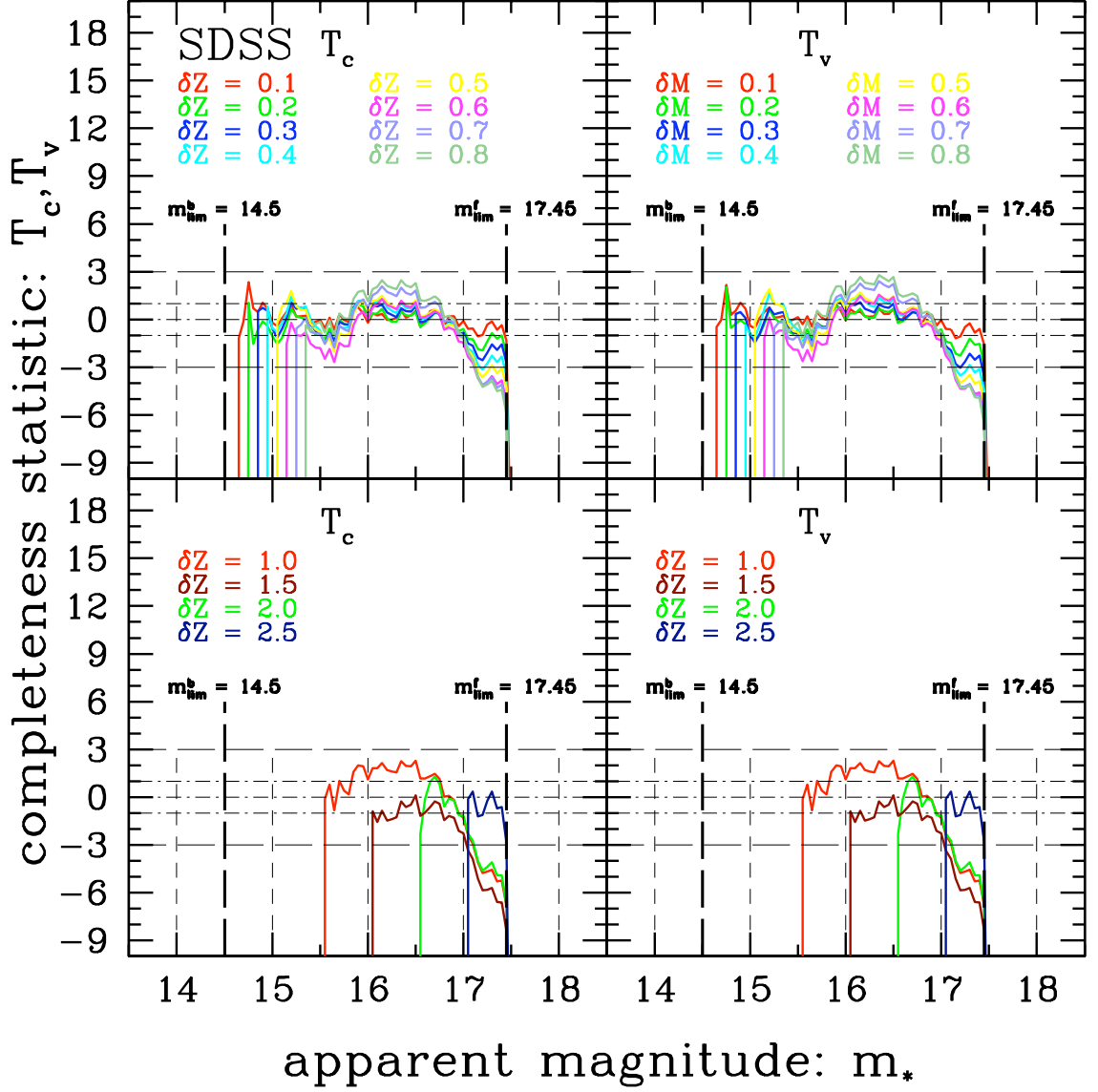


Figure 7.8: T_c and T_v results for SDSS applying the JTH generalisation for values in the range of $0.1 < \delta Z, \delta M < 2.0$. In a similar fashion to that of the 2dFGRS, all four panels exhibit a variation in the resulting trial m_{lim} as both δZ and δM increase, such that for the range $0.4 < \delta Z, \delta M < 2.0$, the resulting trial m_{lim} ranges from $17.15 \lesssim m_{\text{lim}} \lesssim 17.45$ respectively.

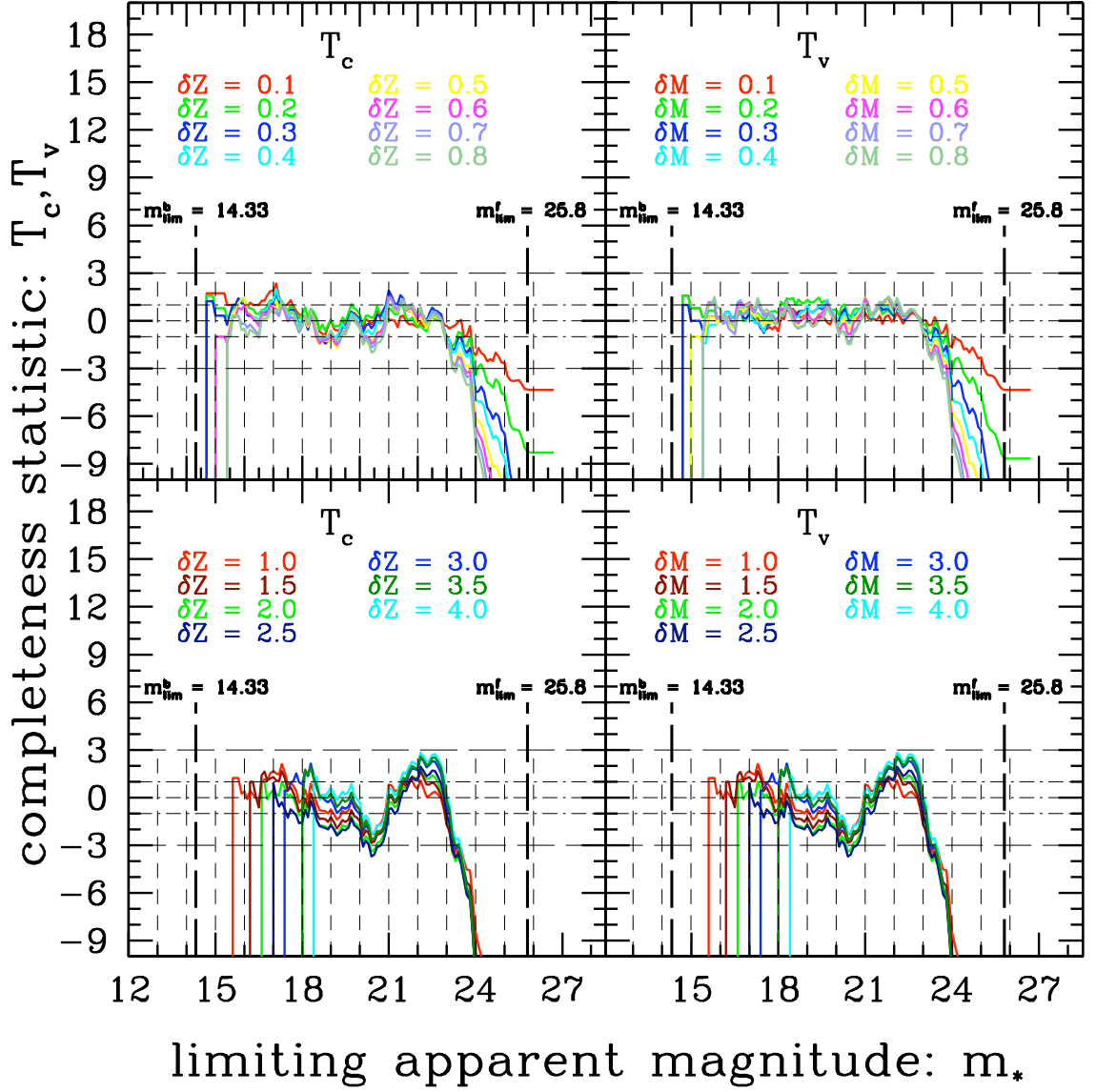


Figure 7.9: T_c and T_v results for the CCLQG data between $0.1 < \delta M, \delta Z < 4.0$. The top panels show T_c and T_v in the range $0.1 < \delta M, \delta Z < 0.8$ where we observe a range in possible limiting magnitudes from $m_* = 25.1$ mag to $m_* = 23.7$ mag as δM and δZ increase. As we move to $1.0 < \delta M, \delta Z < 4.0$, the values for the limiting magnitudes seem to converge to a value of $m_* \sim 23.6$ mag, the same as the when we applied the R01 statistic in the left-hand Figure 6.4.

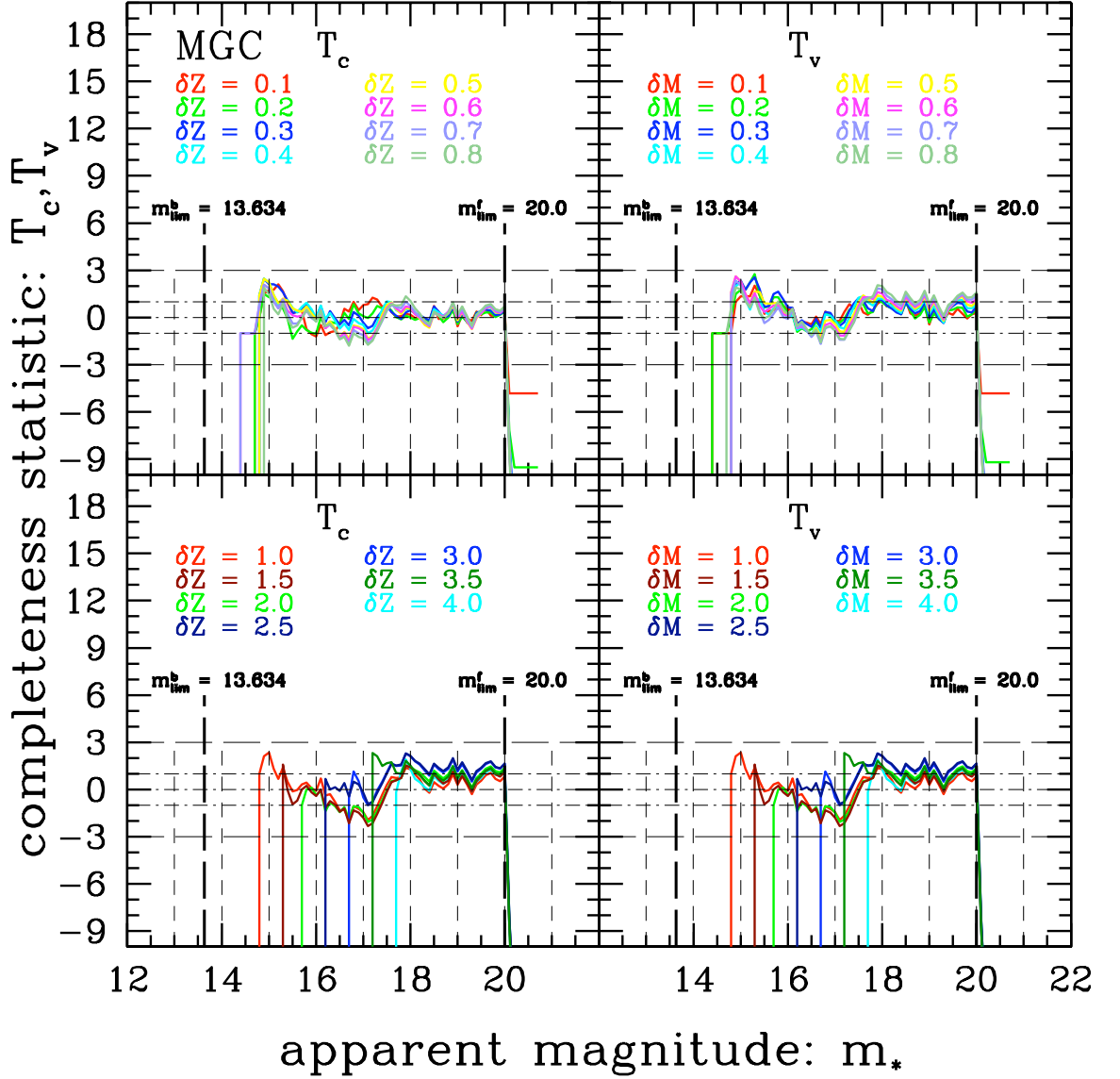


Figure 7.10: T_c and T_v results for the MGC applying the JTH generalisation for values in the range of $0.1 < \delta Z, \delta M < 4.0$. The main difference from that of 2dFGRS and SDSS-Early Types is that no matter how much you increase δZ and δM , the resulting magnitude limit is the same i.e. $m_{\text{lim}} = 20.0$. As we have discussed previously, this data-set seems to be well described by a faint limit only, thus for large δZ and $\delta M \approx 1.0$, T_c and T_v resemble the same characteristics as when the R01 method was applied.

7.2 Optimising the ζ and τ Estimators

In this section we present work that is currently ongoing that concerns optimisation of the JTH estimators, ζ and τ in a way that will overcome the issue of shot-noise detailed in § 7.1.1. Although there are several ways in which we can improve the ζ and τ estimators, we consider two approaches that can be summarised as follows:

1. Maintain a constant number of galaxies, N_{gal} , in the $S_1 \cup S_2$ region for ζ and $S_3 \cup S_4$ region for τ .
2. Calculate T_c and T_v based on setting minimum thresholds respective signal-to-noise ratios of the random variables, $(\zeta/\delta\zeta)$ and $(\tau/\delta\tau)$.

As we will show, the second scenario potentially provides us with a much more robust optimiser than the first. However, the inception of the signal-to-noise (s/n) approach came directly from the results from our initial trials for having a fixed number of galaxies in the ζ and τ regions which is shown to be a rudimentary s/n approach.

7.2.1 Approach #1: Fixing N_{gal}

Firstly, we denote the number of galaxies in $S_1 \cup S_2$ (n_i region) for ζ and $S_3 \cup S_4$ (t_i region) for τ by, N_{gal} . When estimating ζ and τ we can, for every galaxy located at (M_i, Z_i) , maintain a fixed number of galaxies, N_{gal} , counted in the respective regions, n_i and t_i . Therefore, instead of maintaining a constant width for δZ and δM as in the JTH method, we modify JTH (and if we wished, R01 as well) to allow these quantities to vary. However, as we will demonstrate, this approach introduces a trade off between accurately determining the true faint apparent magnitude limit m_{lim}^f and the testing completeness for the whole possible ranges in m_* for a given survey, particularly at brighter m_* 's where there are inherently fewer galaxies observed.

In Figure 7.11 we have applied this procedure to the R01 method to MGC, where the absence of an imposed bright limit provides a simple testing ground for this form of optimisation. We shall discuss the effect of a bright limit on optimising our estimators in more detail in the following sub section. The left- and right-hand panels show the resulting respective T_c and T_v curves for a trial range, $10 \leq N_{\text{gal}} \leq 500$. For both T_c and T_v with $N_{\text{gal}} \lesssim 100$ we observe similar behaviour to that of the JTH method for small fixed widths of δZ and δM where shot-noise dominates, resulting in flat-lining beyond the survey data m_{lim} . Therefore, if we wish to accurately determine the true completeness limit of the survey, where our test statistics systematically drop below

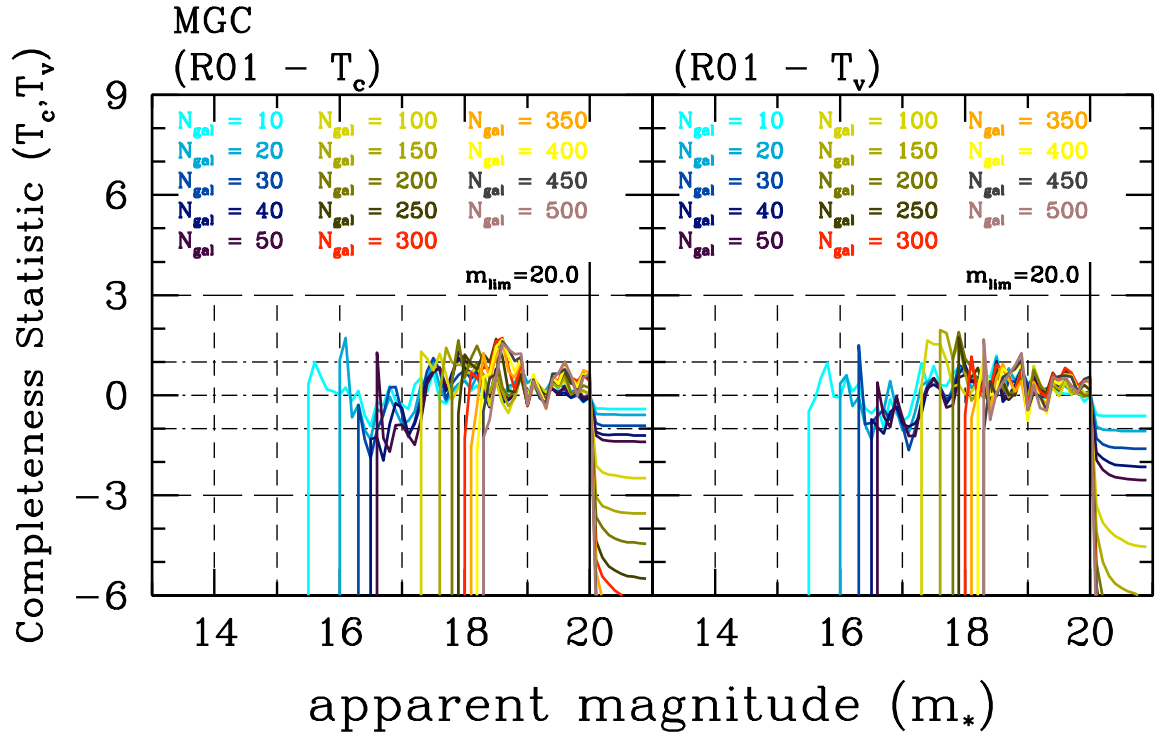


Figure 7.11: MGC T_c and T_v results applying a modified version of the R01 method where the number of galaxies, N_{gal} , in the respective $S_1 \cup S_2$ and $S_3 \cup S_4$ regions is kept constant. In this example we have applied successive trial values of N_{gal} from 10 to 500 galaxies and present the resulting T_c and T_v curves. For T_c in the left panel we observe that up to $N_{\text{gal}} \sim 100$, the statistic is shot-noise dominated. Similarly, for T_v , shot-noise dominates the statistic up to $N_{\text{gal}} \sim 50$. As we move to greater values of N_{gal} up to 500 galaxies, we observe both statistics systematically drop below -3σ . Therefore in order to accurately determine the true apparent magnitude limit we must have a relatively large number of galaxies in the ζ and τ regions. However, as we can see in both panels, this affects the sampling at brighter m_* 's where there are not enough galaxies to be included in the overall T_c and T_v calculations.

-3σ , we require to increase the number of galaxies in N_{gal} . In both panels in Figure 7.11 we can see that for sufficiently large $N_{\text{gal}} \gtrsim 100$ both T_c and T_v begin to drop below -3σ as required by the method. However, this comes at a price. If we again look closely at both panels on Figure 7.11 we also observe that as we increase the number of galaxies, N_{gal} , required to estimate ζ and τ , our ability to ascertain the completeness limit at brighter m_* 's diminishes. This simply implies that as N_{gal} increases there are not enough observable brighter galaxies available to estimate ζ and τ . This results in these galaxies being dropped from the overall T_c and T_v calculation.

In Figure 7.12 we look more closely at how δZ and δM vary with m_* for two of the trial fixed number of galaxies used in MGC. For T_c in the top panels, $N_{\text{gal}} = 10$ and 150,

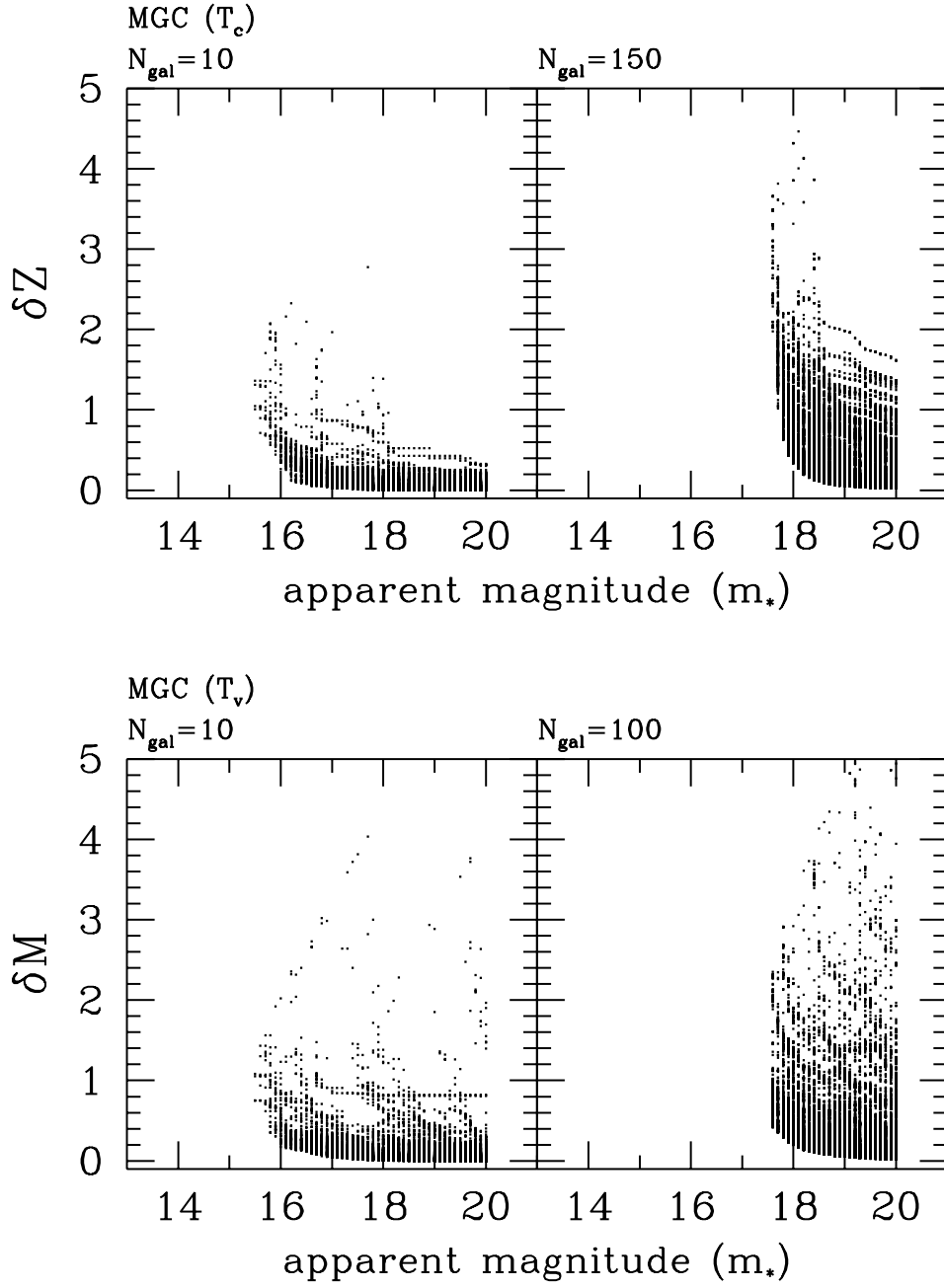


Figure 7.12: Resulting δZ and δM distribution for the respective MGC T_c and T_v statistics. The top panels show how δZ (for T_c) varies at every trial apparent magnitude, m_* , for $N_{gal} = 10$ (left) and $N_{gal} = 150$ (right). For $N_{gal} = 10$ this represents the shot-noise dominated statistic shown in the left panel of Figure 7.11. As we can see this results in the majority of δZ values being distributed between $0 < \delta Z \lesssim 0.2$. Consequently, for $N_{gal} = 150$ (top right panel), the overall δZ increases to $0 < \delta Z \lesssim 2.0$. We observe a similar trend for T_v in the bottom panels.

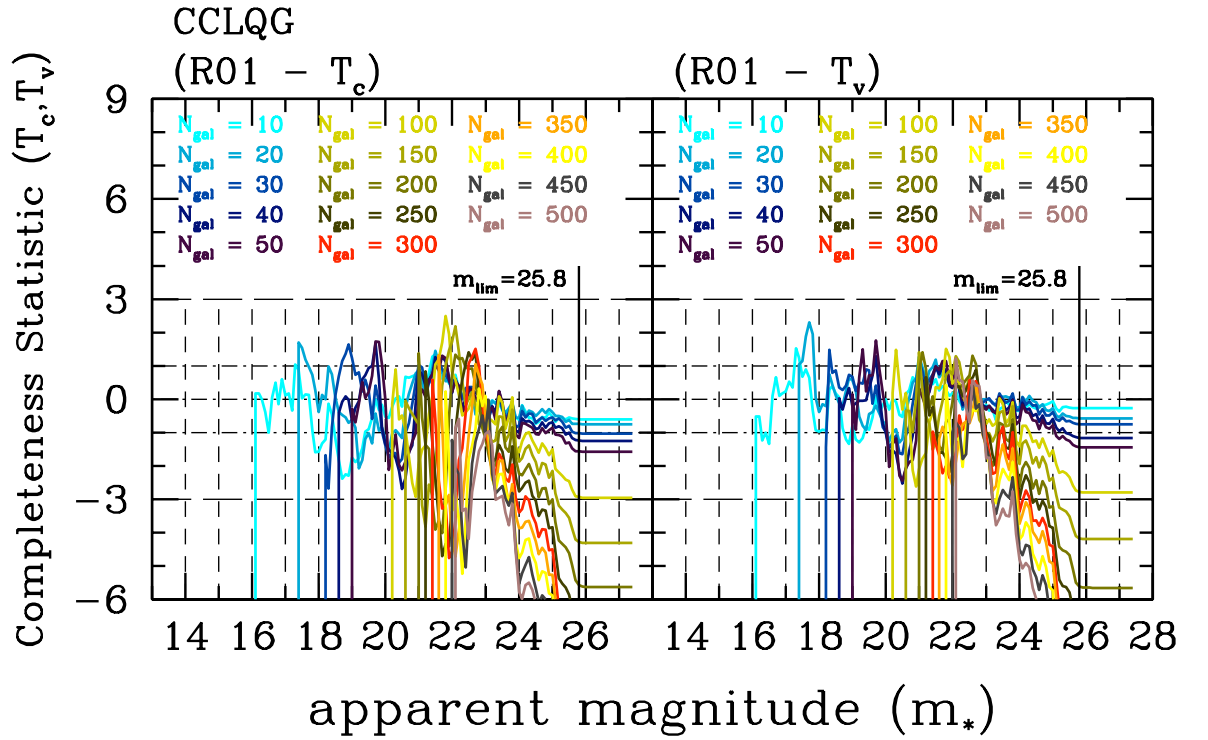


Figure 7.13: CCLQG T_c and T_v results applying a modified version of the R01 method where the number of galaxies, N_{gal} , in the respective $S_1 \cup S_2$ (left panel) and $S_3 \cup S_4$ (right panel) regions is kept constant. The CCLQG has a comparative number of galaxies in the total data-set as MGC and therefore we vary N_{gal} by the same amount as in Figure 7.11. For this survey, we observe the same behaviour as in MGC where, for small values of $N_{\text{gal}} \lesssim 100$, both T_c and T_v are noise dominated and therefore is impossible to accurately determine the true apparent magnitude limit. For each successive N_{gal} up to 500, we observe the same variation in the indication of the true m_{lim} as we have already seen in Figure 7.9 (on page 124). However, by modifying the *ROBUST* method in this way allows us, in principle to determined a fixed error on T_c and T_v for each fixed value of N_{gal} .

and for T_v in the bottom panels $N_{\text{gal}} = 10$ and 100. Each point in the plots represent the final δZ or δM value for the resulting respective ζ or τ calculation. As we would expect, for the small value of $N_{\text{gal}} = 10$, the majority of the respective δZ and δM distributions are approximately in the range $0 < \delta Z \lesssim 0.2$ and $0 < \delta M \lesssim 0.6$. For the larger N_{gal} values that result in the test statistics dropping systematically below -3σ (as shown in Figure 7.11), the overall δZ and δM have to increase in order that they include the required number of galaxies. For $T_c(N_{\text{gal}} = 100)$ we observe the majority of δZ in the range $0 < \delta Z \lesssim 2.2$, and for $T_v(N_{\text{gal}} = 150)$ the majority of δM are in the range $0 < \delta Z \lesssim 1.5$.

We have applied this approach to the CCLQG data-set. The T_c and T_v results in

Figure 7.13 show a similar trend to that observed in MGC. Here, we have applied the same range of N_{gal} for both ζ and τ estimators as in MGC and we can see that as the trial N_{gal} values increase beyond ≈ 100 , both T_c and T_v drop below -3σ indicating where the true m_{lim} is located. However, as we have already observed with this survey, when applying the JTH method, as N_{gal} increases the indicated magnitude completeness limit also varies. However, the scope for this approach would in principle allow us to determine an error for each T_c and T_v point. It should also be reinforced that although setting the number of galaxies, N_{gal} , to be included in each estimation of ζ and τ to increasingly larger values, increases the amount of signal over the noise, it severely limits our ability to test the completeness limits of the whole data-set. For example, if we look at the the $T_c(N_{\text{gal}} = 500)$ curve for MGC in Figure 7.11, it is obvious to see that the test statistic does not begin to register a result until $m_* \approx 18.4$ mag which equates to approximate range of 4.5 mags from the brightest galaxy in the data-set. Therefore, if we are to retain the essence of the Rauzy method whereby one can assess the completeness limit over the whole apparent magnitude range of a given survey, then we require a more sophisticated approach that allows us to use as much of the data as possible without being noise dominated and being able to determine an error estimate on the T_c and T_v statistics.

7.2.2 Approach #2: Measuring the signal-to-noise

Our proposed second solution to this problem lies in estimating the random variables, ζ and τ based on the calculated signal-to-noise (s/n). As we will demonstrate, the s/n -ratio can tell us the minimum number of galaxies we require to accurately estimate ζ and τ . Therefore, T_c and T_v can be calculated based either on a minimum or constant s/n level.

For the T_c statistic we have derived an expression for the s/n -ratio where ζ represents the signal and $\delta\zeta$ represents the noise. In this case Equation 7.1 then becomes

$$\delta\zeta = \frac{\delta r_i(n_i + 1) - r_i\delta(n_i + 1)}{(n_i + 1)^2}. \quad (7.3)$$

To take into account the cross-terms we square Equation 7.3,

$$(\delta\zeta)^2 = \frac{\delta r_i^2}{(n_i + 1)^2} + \frac{\zeta^2[\delta(n_i + 1)]^2}{(n_i + 1)^2} - \frac{2\zeta\delta r_i[\delta(n_i + 1)]}{(n_i + 1)^2}, \quad (7.4)$$

$$\frac{\zeta^2}{(\delta\zeta)^2} = \frac{r_i^2}{\delta r_i^2} + \frac{(n_i + 1)^2}{[\delta(n_i + 1)]^2} - \frac{r_i(n_i + 1)}{2\delta r_i[\delta(n_i + 1)]}. \quad (7.5)$$

Therefore, the s/n ratio is for ζ is given by,

$$\frac{\zeta}{(\delta\zeta)} = \left[\frac{r_i^2}{\delta r_i^2} + \frac{(n_i + 1)^2}{[\delta(n_i + 1)]^2} - \frac{r_i(n_i + 1)}{2\delta r_i[\delta(n_i + 1)]} \right]^{1/2}. \quad (7.6)$$

By applying a similar approach for T_v we can obtain a similar expression for the signal-to-noise for estimating τ . Starting from Equation 7.2 we can show that,

$$\frac{\tau}{(\delta\tau)} = \left[\frac{q_i^2}{\delta q_i^2} + \frac{(t_i + 1)^2}{[\delta(t_i + 1)]^2} - \frac{q_i(t_i + 1)}{2\delta q_i[\delta(t_i + 1)]} \right]^{1/2}. \quad (7.7)$$

From a computational point of view, the incorporation of this approach defined by Equations 7.6 and 7.7 into the current JTH method would be tricky and computationally intensive as we have illustrated with the SDSS M - Z distribution in Figure 7.14. Recalling the current construct of the JTH method, for a fixed value of δZ we can determine the quantity, M_{lim}^b which allows us to exclude galaxies in the region we have denoted as, S_T , and define the rectangular region $S1 \cup S2$ which remains between m_* and M_{lim}^b . However, if, for example, we wish to base our estimation of ζ on a fixed s/n -ratio level then we have to allow δZ to grow in size until the appropriate s/n value has been reached. This implies that for every incremental increase in δZ (left-hand panel of Figure 7.14) we have to continually re-define the quantity M_{lim}^b in order to remove any galaxies that may be present in the S_T region. This only becomes a significant issue for larger δZ , as illustrated on the right-hand panel of Figure 7.14. It is this extra process that would significantly hamper the efficiency of the current coding of T_c (and T_v). Another important issue we have to consider is how the growth of δZ affects the size of the S_1 and S_2 and the resulting s/n . In the left-hand panel of Figure 7.14 we observed that as δZ increases in size, the S_1 and S_2 regions also grow, increasing the number of galaxies going in to the ζ calculation. However, as δZ increases to the indicated size on Figure 7.14 as Z_{i+n} , S_1 begins to narrow due to the presence of m_{lim}^b whilst S_2 continues to grow. It follows, therefore, that there would be a point at which δZ becomes so large that S_1 becomes too narrow to sample any galaxies. As we will see, this effect introduces what we can think of as a ‘forbidden’ region on the $\delta Z - m_*$ plane, allowing us to set limiting values on δZ . The exact same issue is also inherent in the τ estimation for increasing δM .

7.2.2.1 The growth of s/n : initial tests

To better understand how we might expect the s/n grow with m_* we initially incorporated Equations 7.6 and 7.7 into the R01 method. For the reasons described above, the

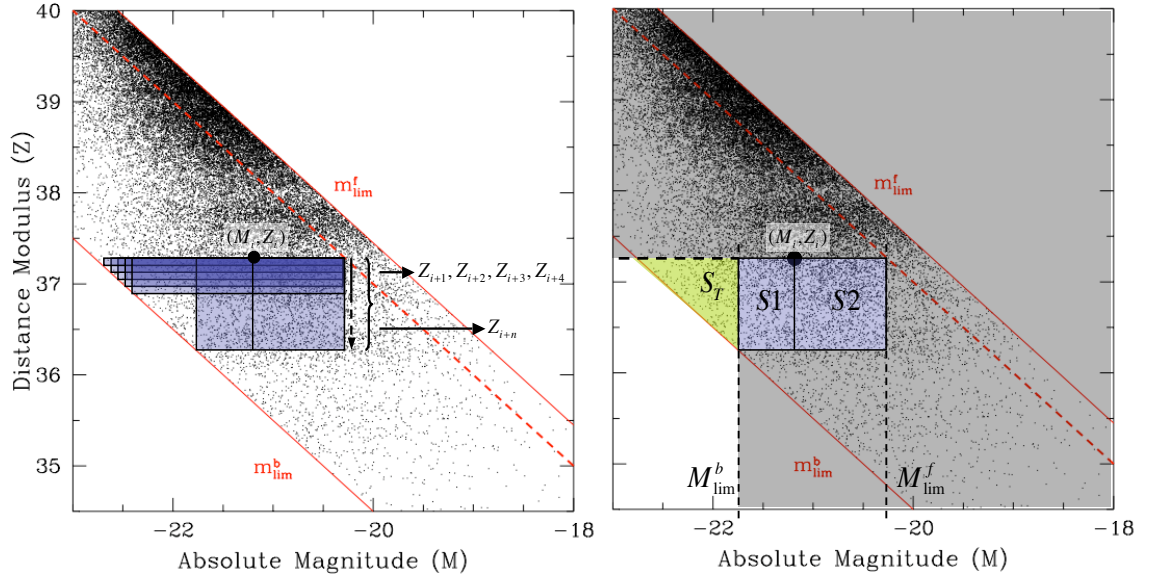


Figure 7.14: Schematic illustrating a procedure to calculate the signal-to-noise for ζ . The left hand plot is a section of the M - Z distribution for SDSS-EDR. Here, we consider the construction for the ζ estimate for a galaxy located at (M_i, Z_i) at a trial magnitude limit $m_* = 17.0$ mag. Since we have already sorted the data-set for Z we begin at Z_i then construct the S_1 and S_2 regions around $Z_i \rightarrow Z_{i+1}$, and calculate the signal-to-noise ratio. We then systematically construct S_1 and S_2 for $Z_i \rightarrow Z_{i+2}, Z_{i+3}, \dots, Z_{i+n}$ until the predetermined signal-to-noise level has been reached. However, constructing S_1 for each incremental shift in Z_{i+n} is tricky. The right hand panel shows the S_1 and S_2 regions for Z_{i+n} . The grey region represents the galaxies that are easily excluded from the calculation by traditional means. However, we have to now continually calculate a new region denoted by S_T since M_{lim}^b is no longer a fixed quantity.

R01 method does not consider a bright limiting magnitude and therefore we can allow δZ and δM grow to there maximum - provided of course we use a suitable data-set. We, therefore, use MGC for this analysis where it has proven throughout our analysis to be a very good calibrator for our estimators.

In Figure 7.15 we have plotted the distribution $\log \zeta / \delta \zeta$ (top panels) and $\log \tau / \delta \tau$ (bottom panels) as a function of δZ and δM respectively. In the left-hand panels in both cases each point represents the total $\log(s/n)$ for the i^{th} galaxy's ζ calculation. The coloured regions indicate the increasing values of m_* . As we would expect, for both ζ and τ as m_* moves to increasingly fainter magnitudes the overall s/n increases with it. This is a result of increasingly more galaxies being included in the ζ and τ calculation. In the right-hand panels for both ζ and τ we have also illustrated the rate of growth of the $\log(s/n)$ for each galaxy at an $m_* = 18.0$ mag, indicated by the blue lines. It is clear to see that by and large in both cases, the $\log(s/n)$ rises sharply within

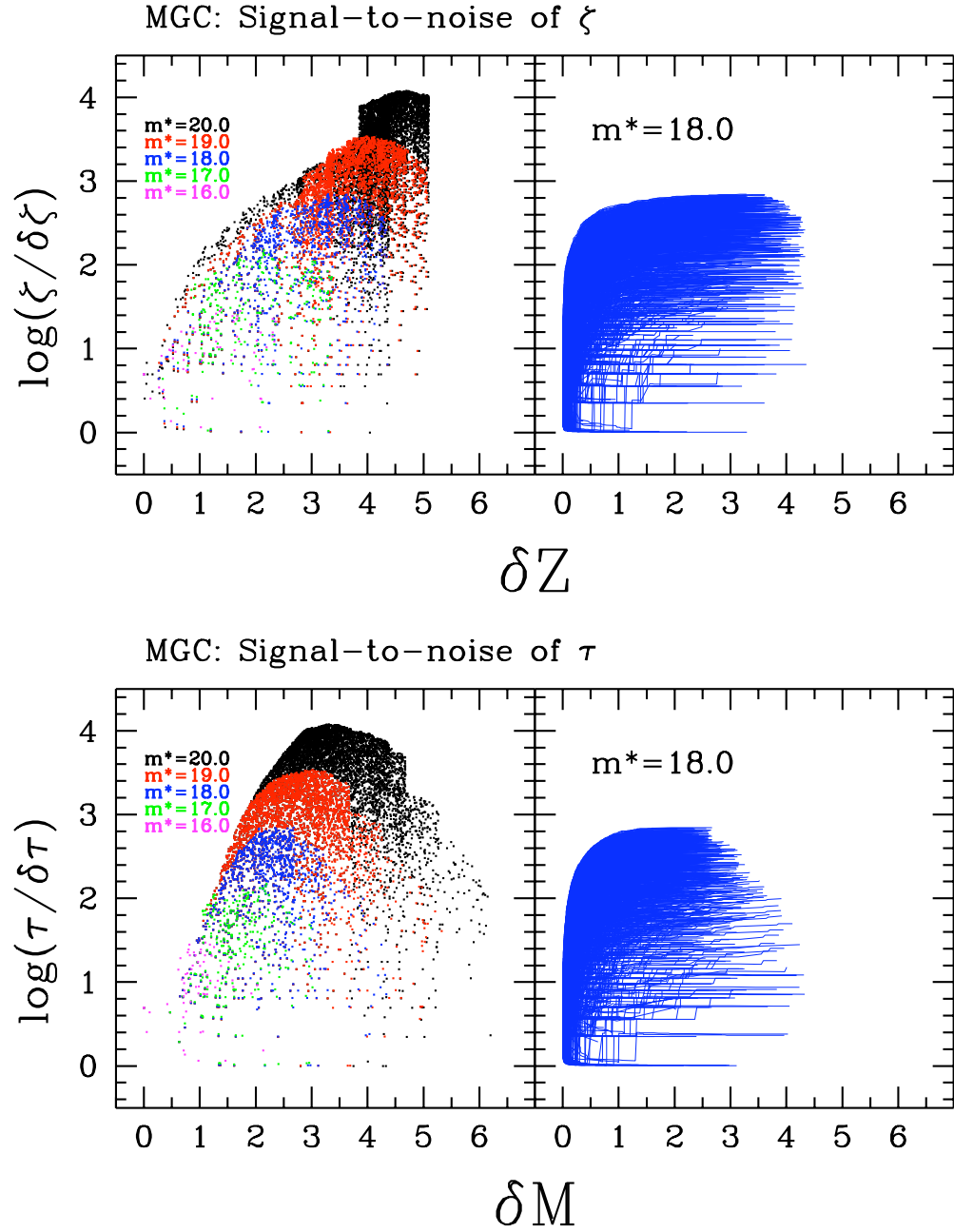


Figure 7.15: MGC survey signal-to-noise (s/n) for the R01 ζ and τ estimators. The top panels show the signal-to-noise for ζ . In the left-hand we have considered the distribution of $\log(s/n)$ [or equivalently $\log(\zeta/\delta\zeta)$] versus varying width δZ . Each point on the distribution represents the total resultant s/n value for the i^{th} galaxy's ζ calculation. Since we are not considering a bright apparent magnitude limit δZ is allowed to grow to its maximum width. The right-hand plot traces growth of s/n for increasing δZ at an $m^* = 18.0$ mag. Therefore, the end point for each trace corresponds to the blue points on the left-hand plot. The bottom panels show the same distributions for the τ estimator.

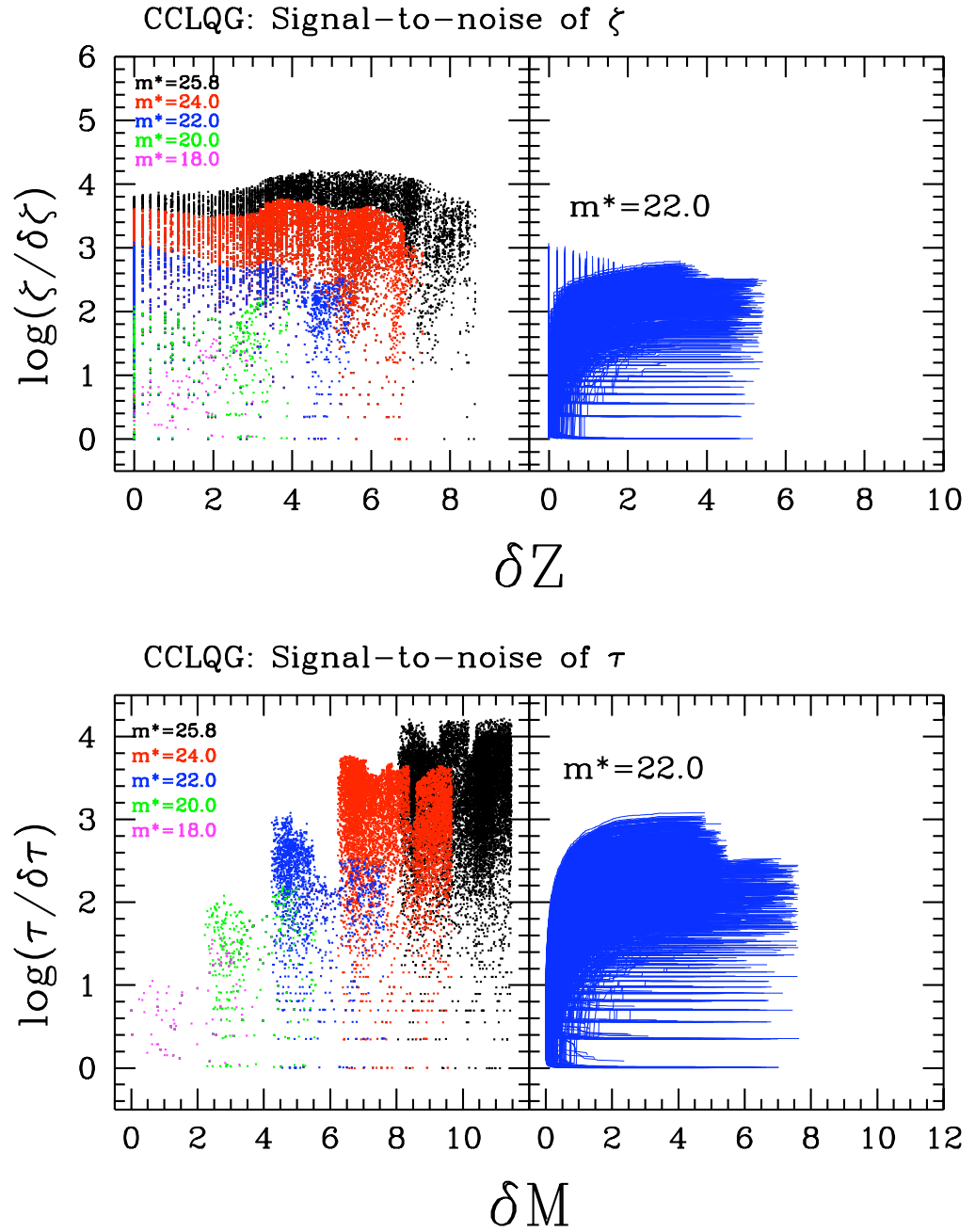


Figure 7.16: CCLQG survey signal-to-noise for the R01 ζ and τ estimators.

a relatively narrow δZ and δM width. This indicates that perhaps a $\delta Z, \delta M \approx 1.0$ is a reasonable width to choose since the bulk of the s/n seems saturated beyond this point, indicating perhaps there is little more relevant information to be added to the calculation for large δZ and δM 's.

We also applied this procedure to the CCLQG data-set (see Figure 7.16) since the

total number of galaxies in the sample is comparable to MGC and we have yet to conclusively prove the need for a bright apparent magnitude limit. We observed in this case the ζ and τ showed very different distributions for their resultant $\log(s/n)$ as m_* . However, in both cases, the overall rising trend of the $\log(s/n)$ is similar to that of MGC which is perhaps better illustrated in the respective right-hand panels of Figure 7.16 which again shows the rate of growth of $\log(s/n)$ for an $m_* = 22.0$ mag. It is interesting to note the discretisation effect of the redshifts manifesting, particularly for smaller δZ , in the top-left panel.

7.2.2.2 From R01 to JTH

The reason for optimising the JTH T_c and T_v statistics was brought about from adverse effects introduced by the presence of a bright limit. However, examining how the s/n varied when applied to R01 method has ultimately lead us to a more concise and efficient way of implementing the s/n to the JTH method. Therefore, for a first order approximation to optimise JTH we plan to determine the s/n for trial values of δZ and δM and for each m_* value, and then interpolate for a threshold s/n value. Consequently, we would have, for each m_* , an averaged δZ and δM width that should, in principle, allow us to utilise as much of the data as possible.

We are presently at the stage of mapping the resulting $\log(s/n)$ as a function of m_* and δZ (or δM) in a slightly more sophisticated way than in Figures 7.15 and 7.16. For our main example, we once again examine MGC shown in Figures 7.17 to 7.20. If we look firstly at Figure 7.17, in both panels we have plotted a 3D surface map of the $\log(s/n)$ varying with m_* and trial values of δZ . For both panels in this figure and indeed in the remaining figures in this chapter we increment m_* , δZ and δM by 0.1. For MGC, the range in δZ that we have considered is $0.1 \leq \delta Z \leq 6.3$, and the range for the trial apparent magnitude is $13.0 \leq m_* \leq 20.0$. The colour gradient represents the changing $\log(s/n)$.

Let us recall for a moment, the distributions in Figures 7.15 and 7.16. In those figures each point represented the total $\log(s/n)$ for each galaxy in the survey sample at a given m_* . In the top panel of Figure 7.17 we have instead determined where the *maximum* (or peak) $\log(s/n)$ occurs for the corresponding m_* and trial δZ width, thereby illustrating the best possible s/n scenario one could achieve for a given sample. However, since the peak s/n is not necessarily the most representative we have also taken the average s/n level shown in the bottom panel of Figure 7.17. The dark blue region observed represents a $\log(s/n) = -1$ which equates to the ‘forbidden’ region

discussed earlier where the δZ or δM width has become too large for the respective S_1 or S_3 region to be sampled. This region is more clearly defined in Figure 7.18. Lastly, we have also superimposed where the maximum $\log(s/n)$ occurs at each m_* for the smallest δZ (and δM across the whole range in δZ (shown in white). In the cases where there are equal values of $\log(s/n)$ we chose the smallest δZ width. If we now look at Figure 7.18 we illustrate further how one might use this line to compute T_c .

The panels in Figure 7.18 are the 2D representation those in Figure 7.17. The black lines in both panels are iso-contours of constant signal-to-noise. If we wish to maximise the efficiency of our code it is our hope that we can use the optimised white line in conjunction with the iso-contours to minimise the width of δZ but retain the optimised s/n . This point is more relevant for the fainter values of m_* where the (M, Z) diagram shows considerably more galaxies beyond $m_* \gtrsim 17.0$ mag. In Figures 7.19 and 7.20 we show the corresponding MGC T_v s/n maps and observe an overall similar trend to that of T_c . For illustrative purposes we have generated the same maps for CCLQG (Figures 7.21, 7.22, 7.23 and 7.24), SDSS (Figures 7.25, 7.26, 7.27 and 7.28) and 2dFGRS (Figures 7.29, 7.30, 7.31 and 7.32).

Future work in this area will be concerned with taking these maps to the next level and applying them to T_c and T_v . Moreover, by exploring the cross-correlations within T_c and T_v we it should be possible to determine their respective covariance which will lead to obtaining an error for the T_c and T_v statistics. The combination of these two areas will provide the user with a statistic that has not only has well defined errors, but also utilises the data in the most efficient and effective way based on the signal-to-noise. The latter point could improve the more established non-parametric estimators.

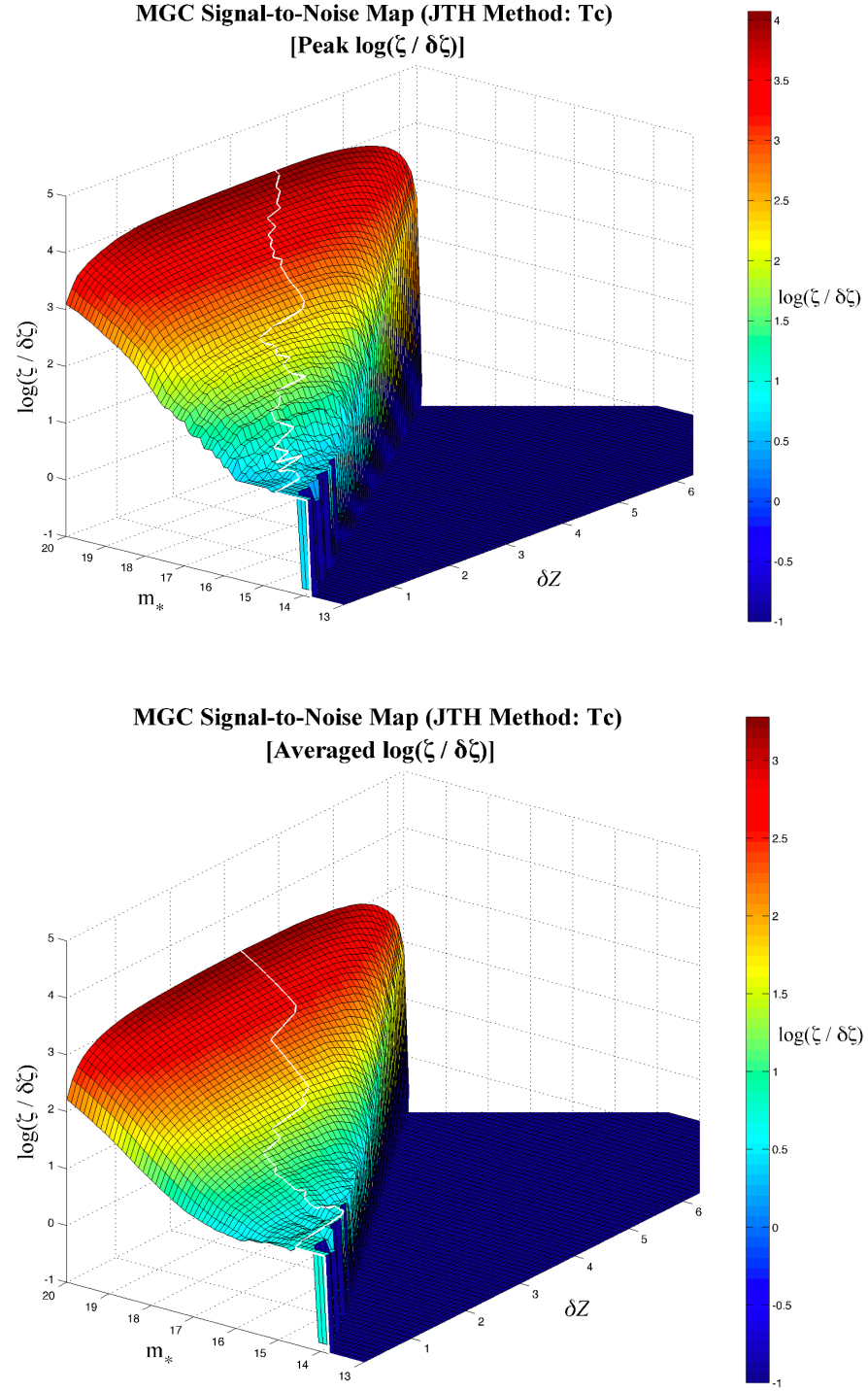


Figure 7.17: 3D representation of the MGC signal-to-noise map derived from the application of the JTH T_c statistic. The map is constructed by applying the JTH T_c statistic in the usual way for successive incremental values of $\delta Z = 0.1$. Then, at each trial apparent magnitude, m_* for a given δZ , we have determined both the peak (top panel) and the average (bottom panel) signal-to-noise. The white line superimposed on both maps charts the maximum signal-to-noise value at each m_* for the smallest δZ allowable over the whole range in δZ . This line therefore, indicates the optimised δZ required to maximise the signal-to-noise.

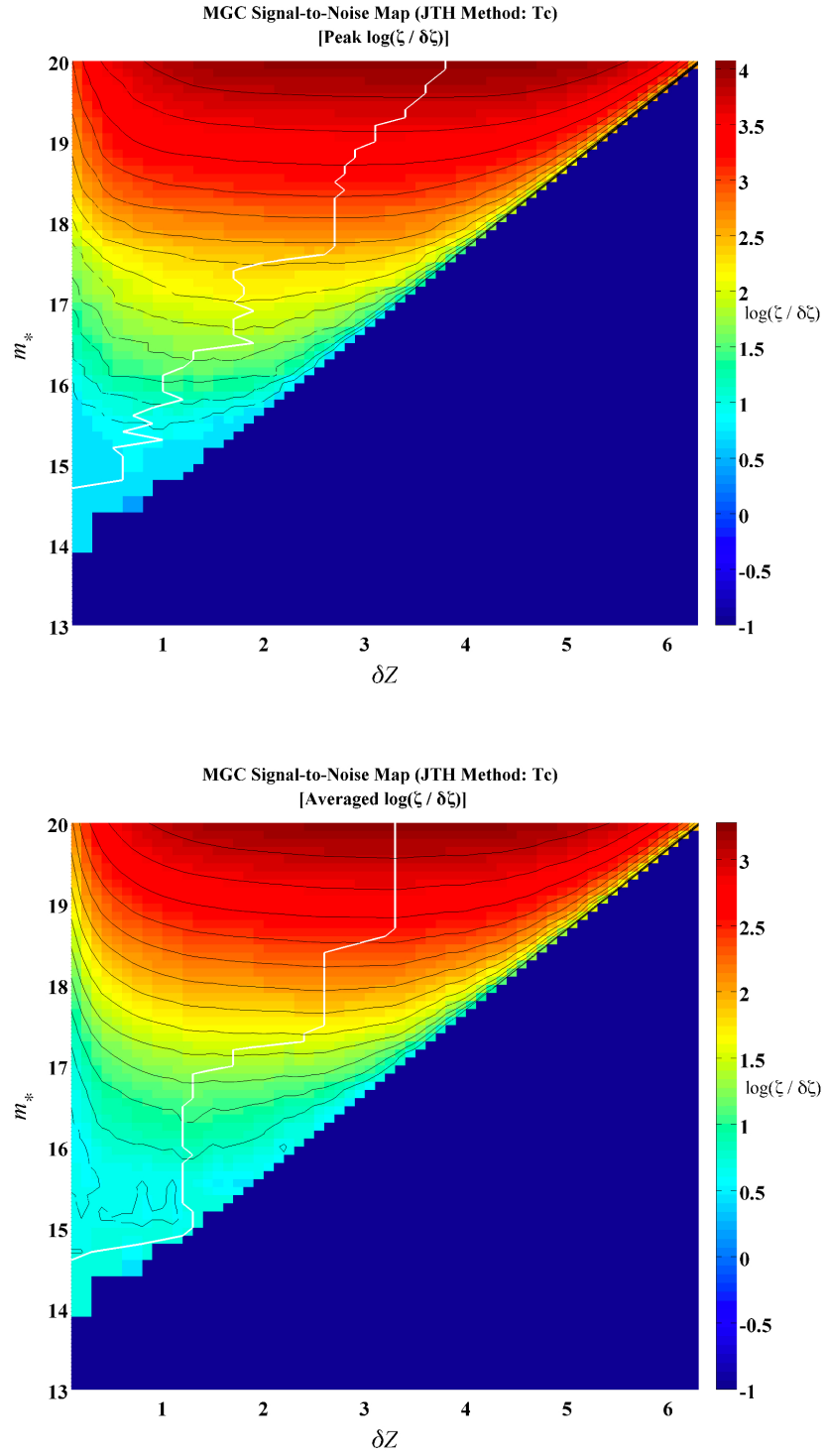


Figure 7.18: 2D representation of the MGC signal-to-noise map derived from the application of the JTH T_c statistic shown in Figure 7.17. The signal-to-noise is shown by the gradated colour. The white line superimposed on both maps indicates the optimised δZ required to maximise the signal-to-noise. The black lines represent iso-contours of constant signal-to-noise. By using the iso-contours in conjunction with the optimised white line we can also maximise the computational efficiency by selecting a smaller δZ value that retains the same s/n level.

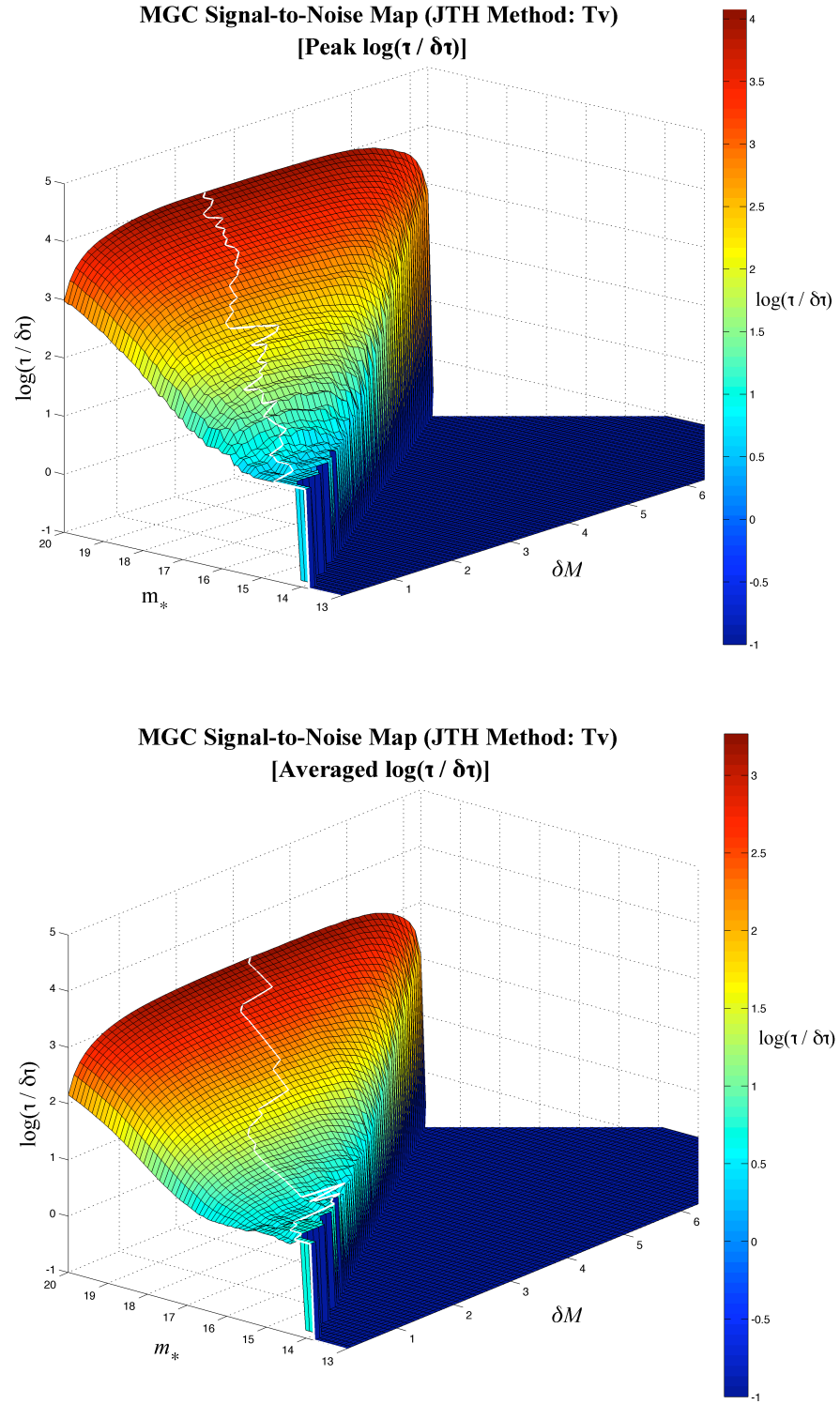


Figure 7.19: 3D representation of the MGC signal-to-noise map derived from the application of the JTH T_v statistic.

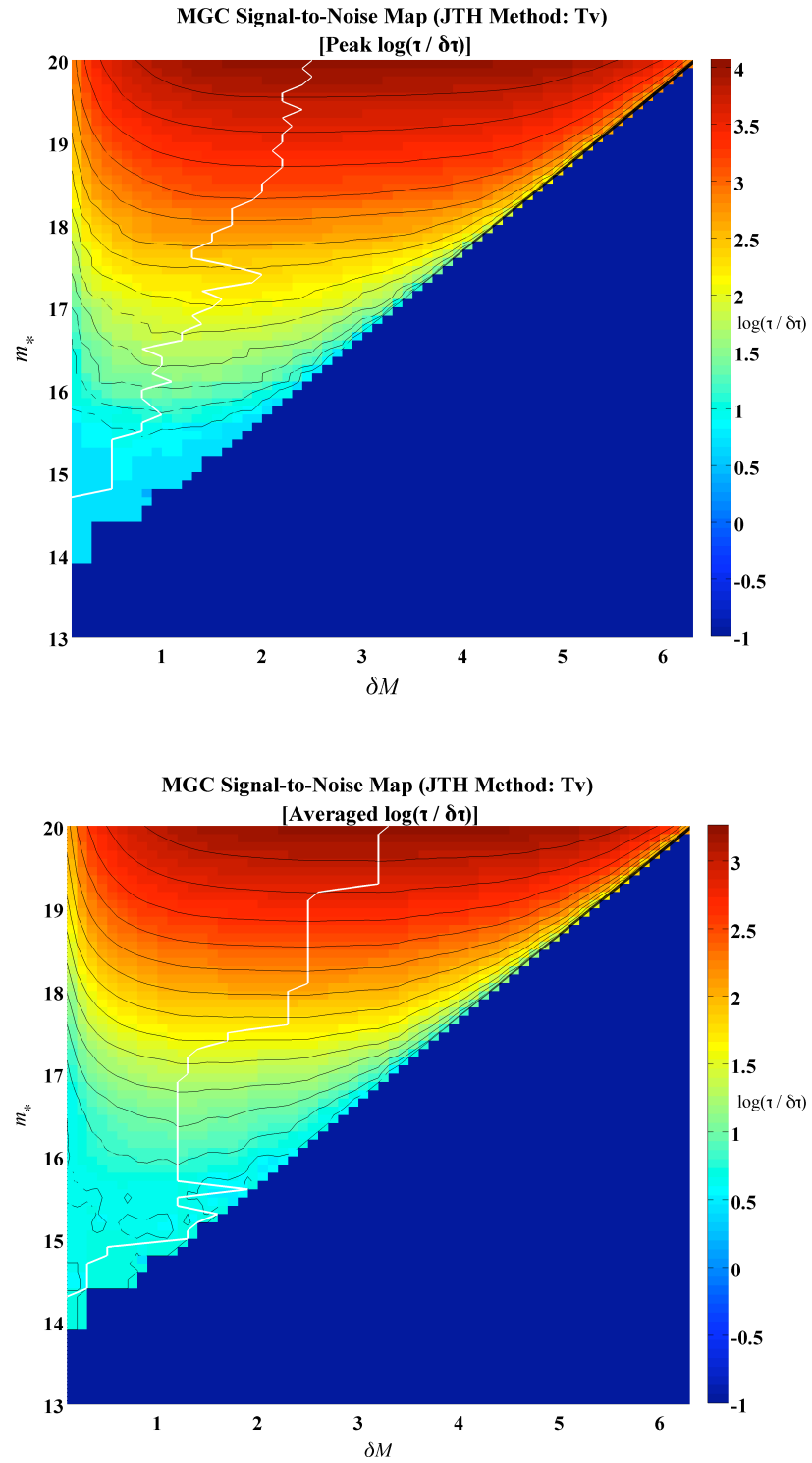


Figure 7.20: 2D representation of the MGC signal-to-noise map derived from the application of the JTH T_v statistic.

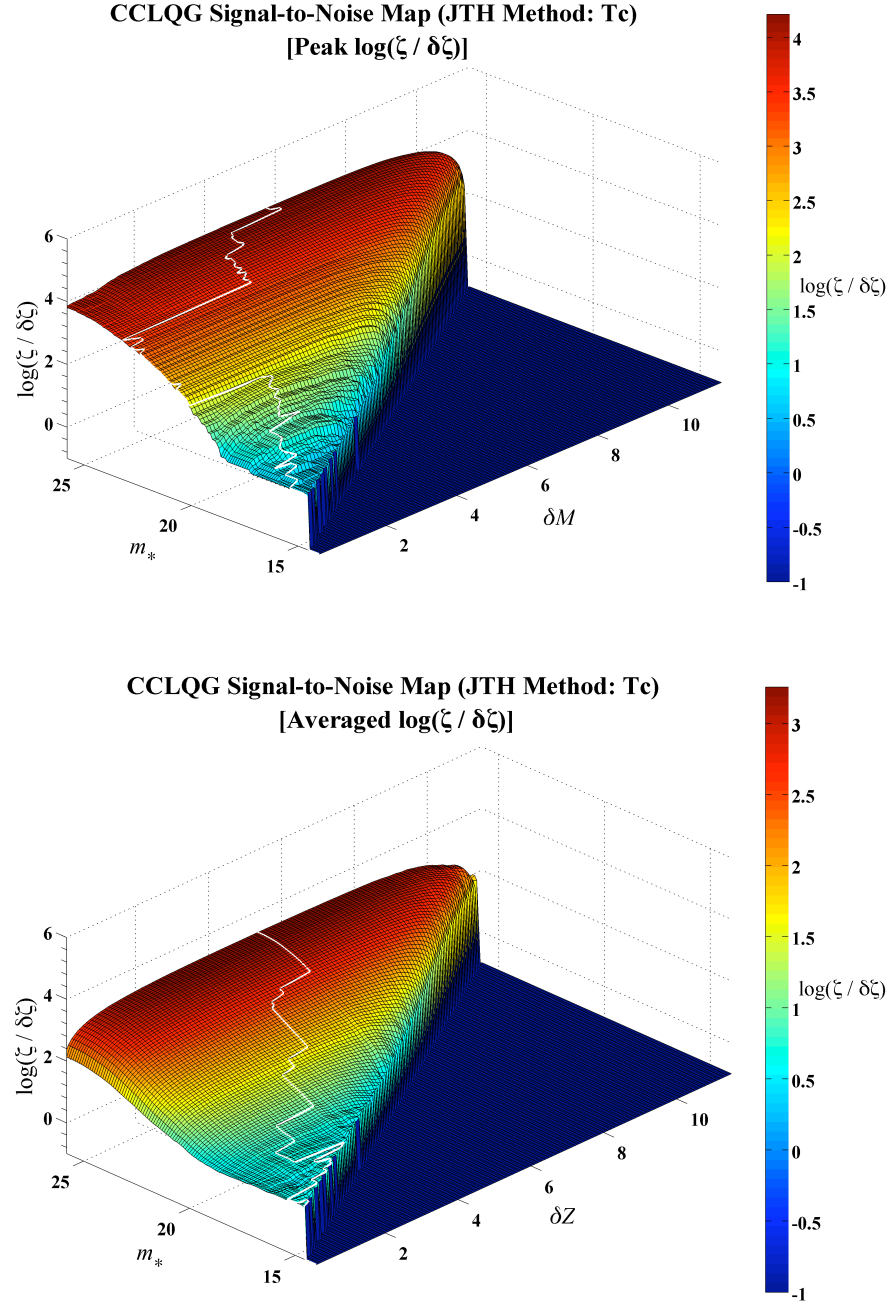


Figure 7.21: 3D representation of the MGC signal-to-noise map derived from the application of the JTH T_c statistic. In the top panel we can see that the white line

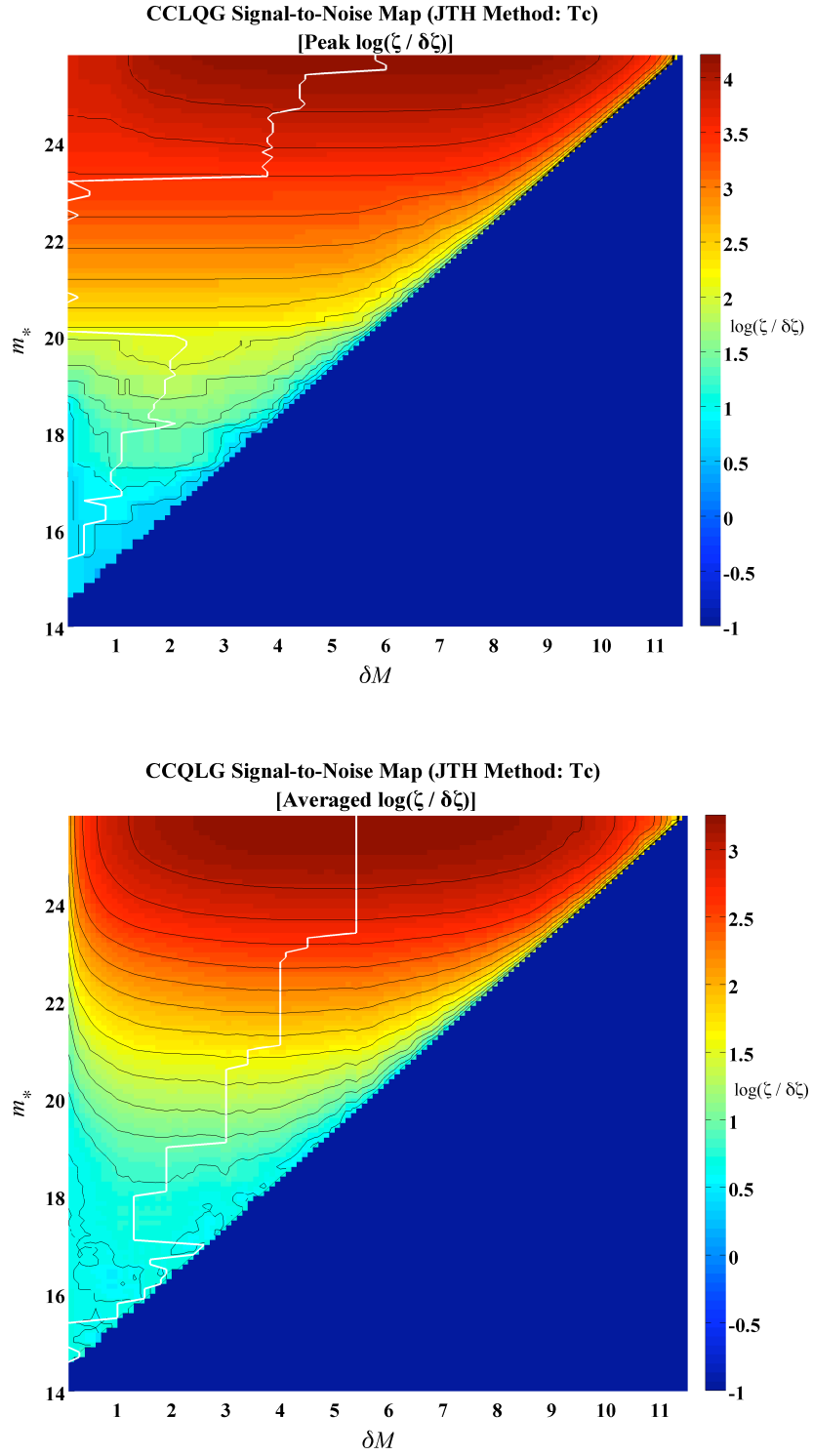


Figure 7.22: 2D representation of the CCLQG signal-to-noise map derived from the application of the JTH T_c statistic.

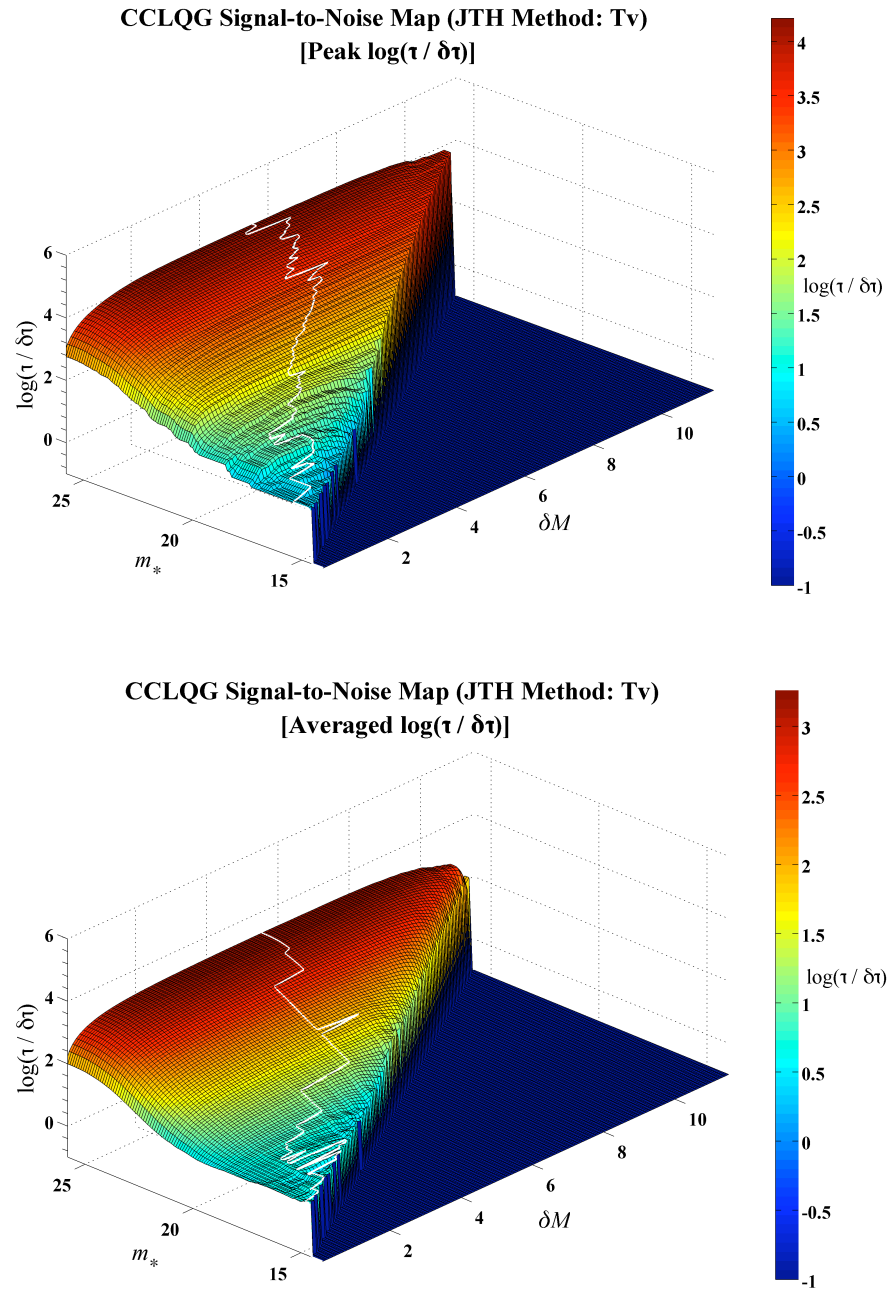


Figure 7.23: 3D representation of the CCLQG signal-to-noise map derived from the application of the JTH T_v statistic.

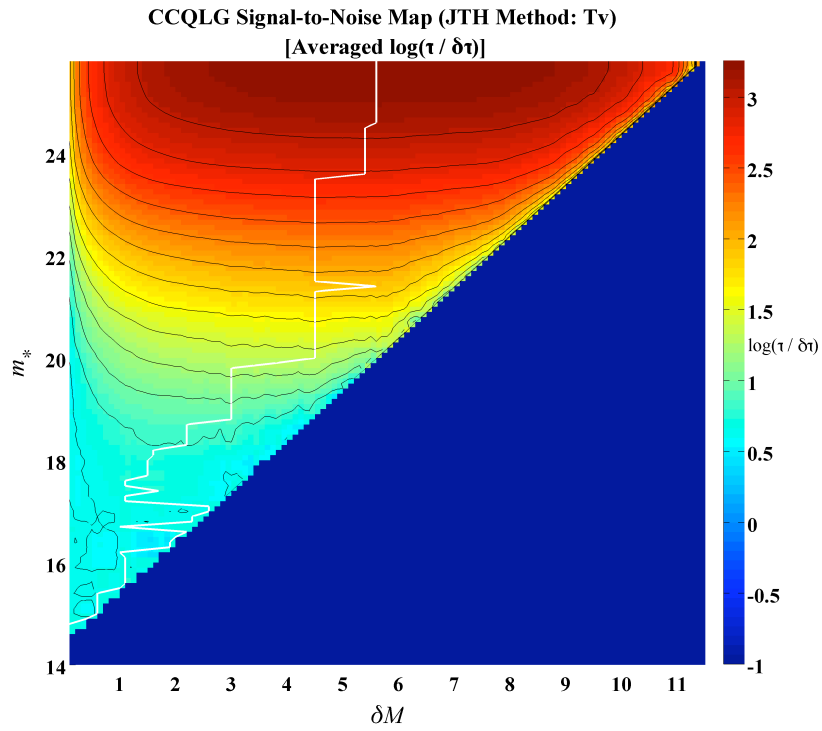
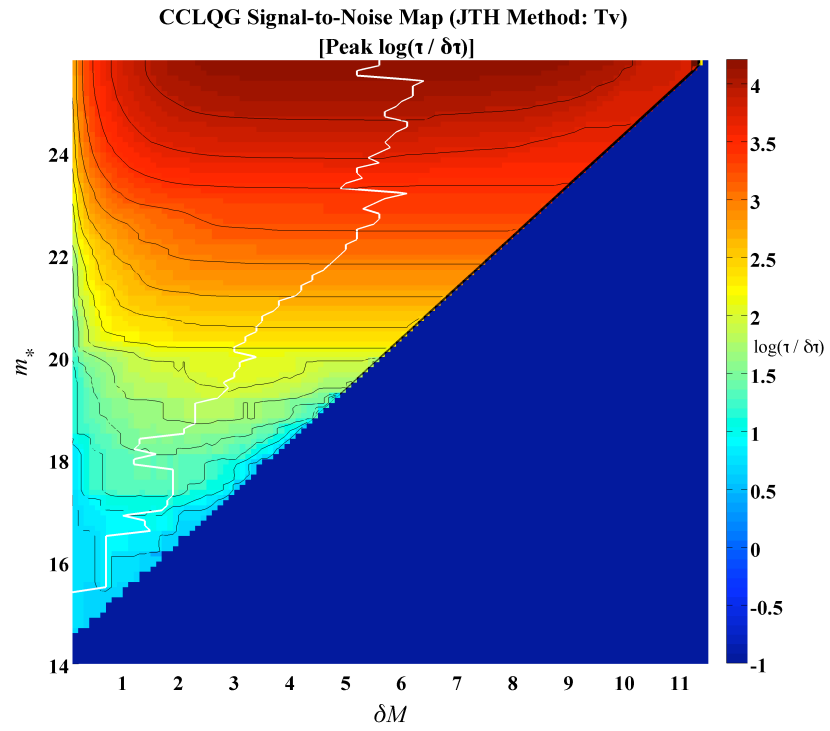


Figure 7.24: 2D representation of the CCLQG signal-to-noise map derived from the application of the JTH T_v statistic.

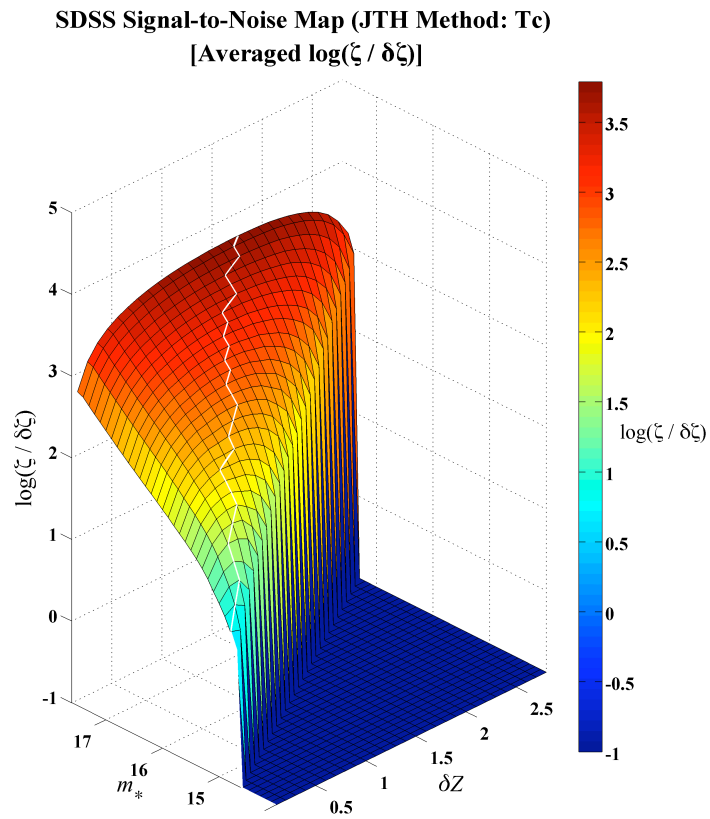
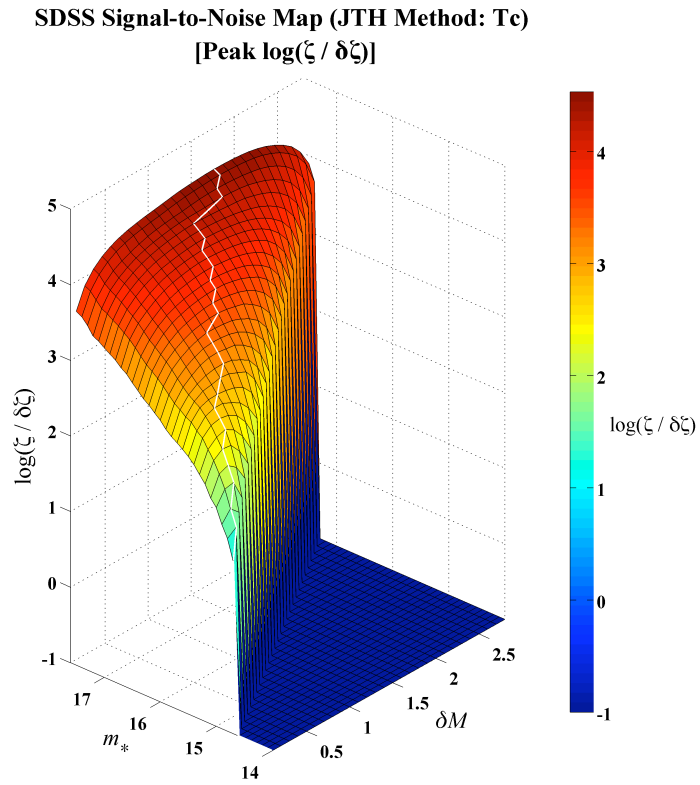


Figure 7.25: 3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_c statistic.

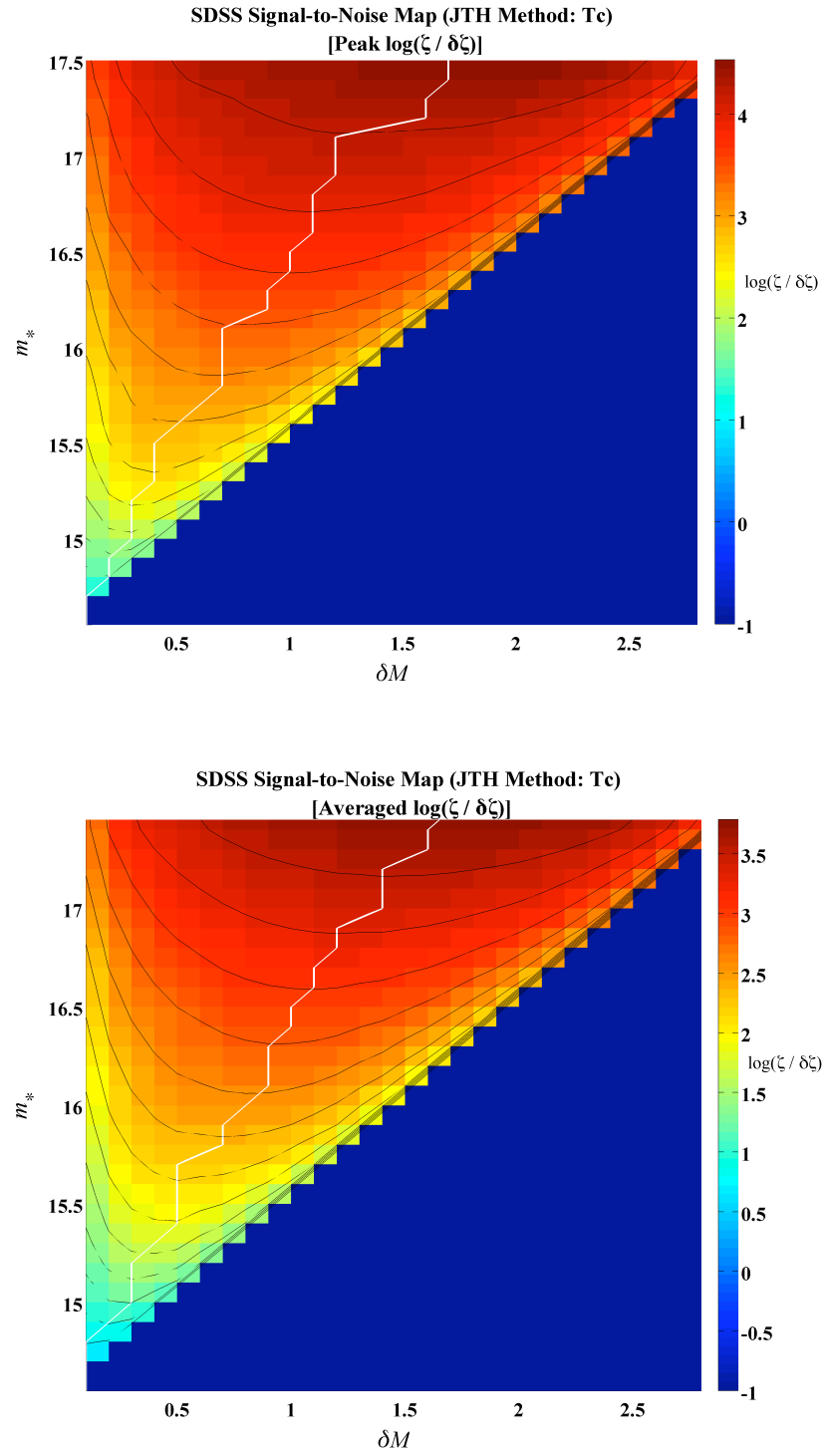


Figure 7.26: 2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_c statistic.

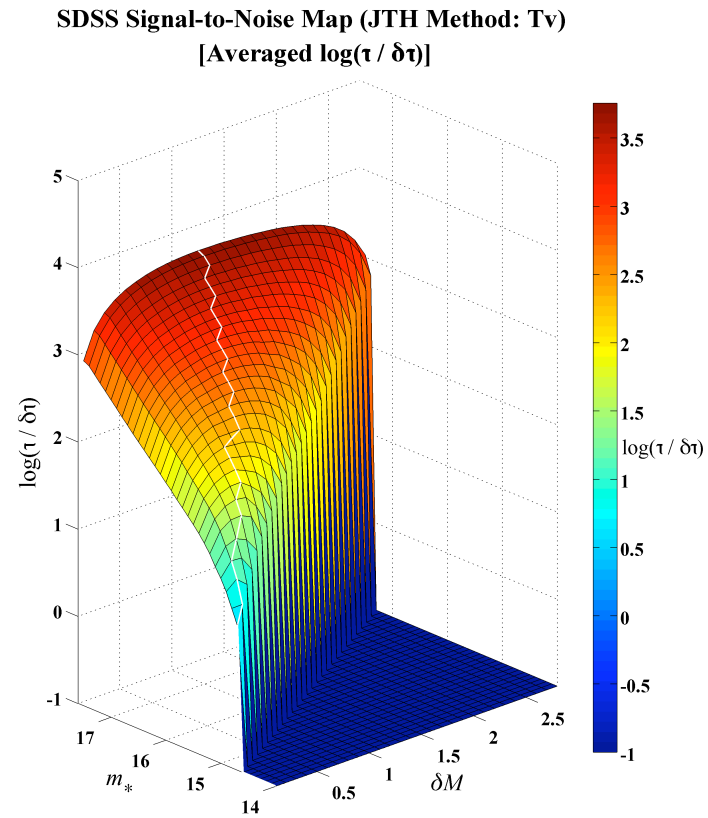
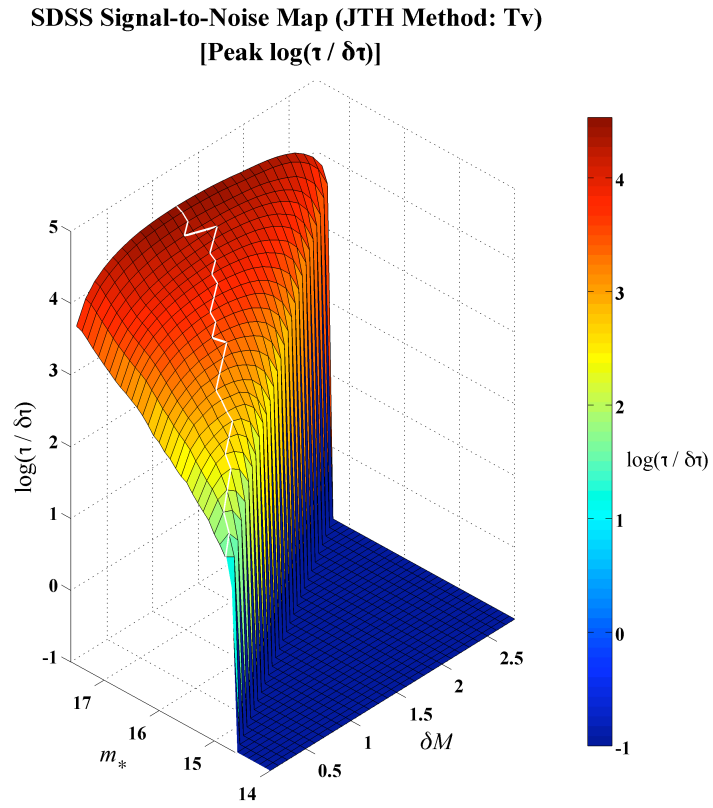


Figure 7.27: 3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.

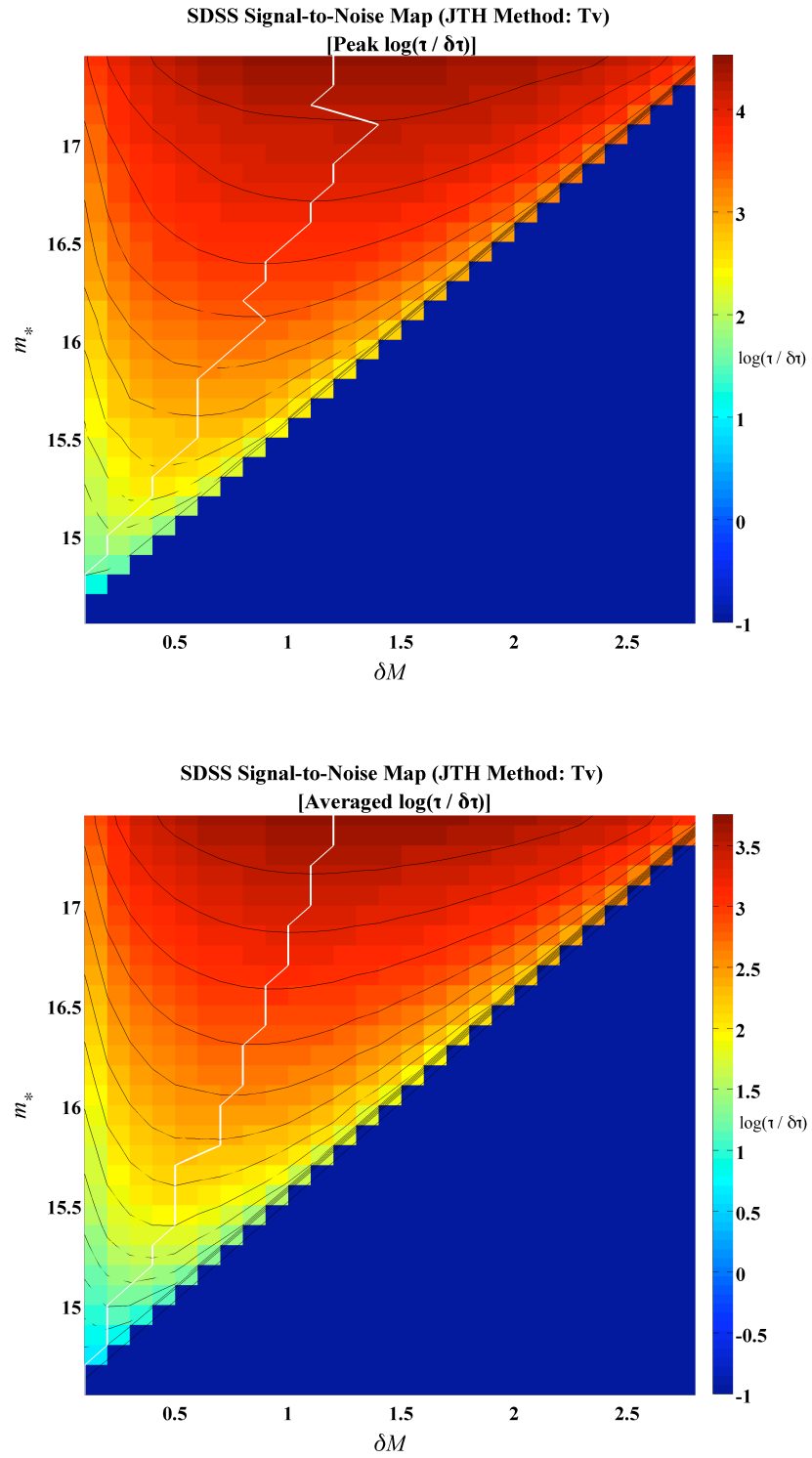


Figure 7.28: 2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.

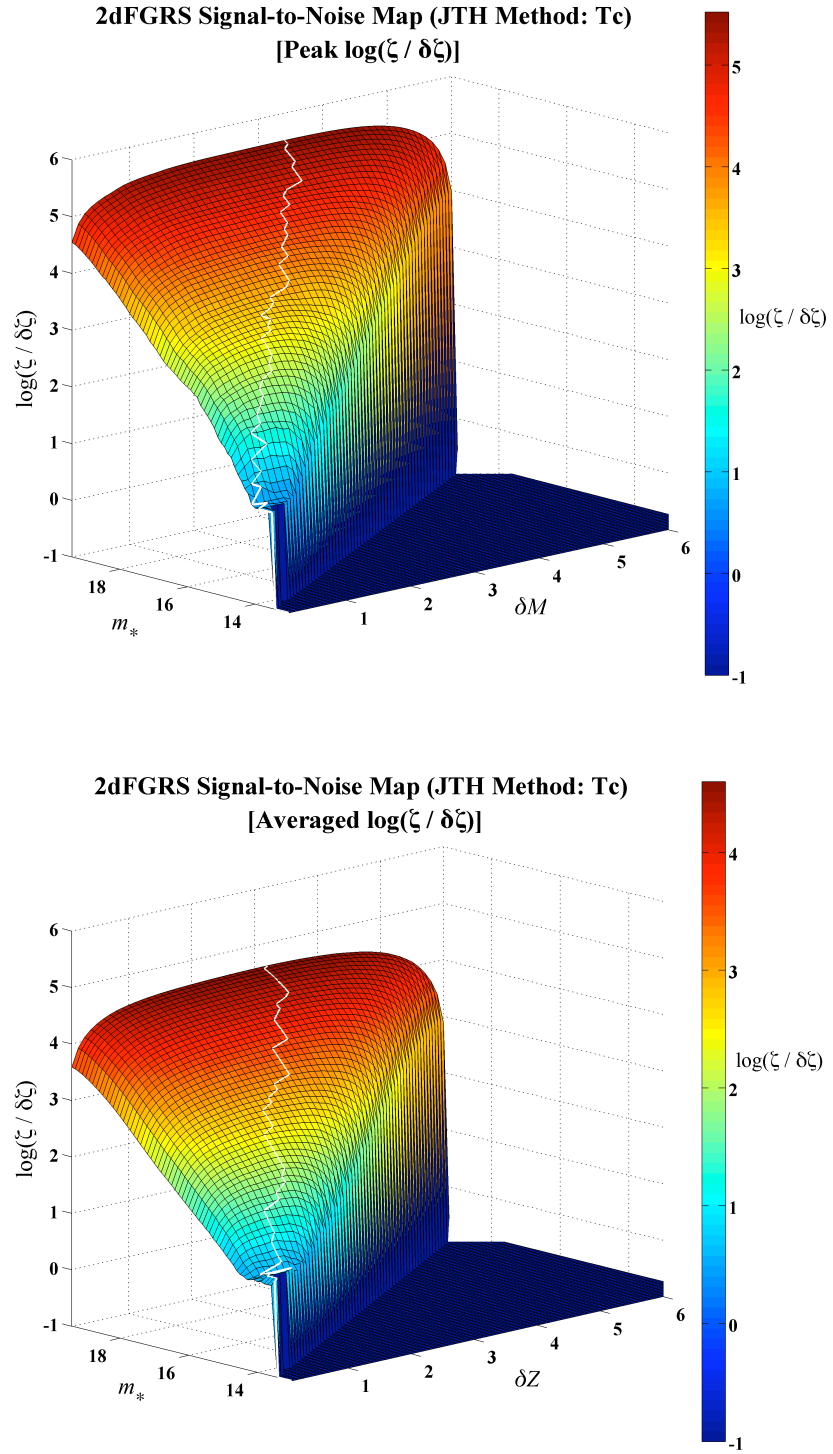


Figure 7.29: 3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.

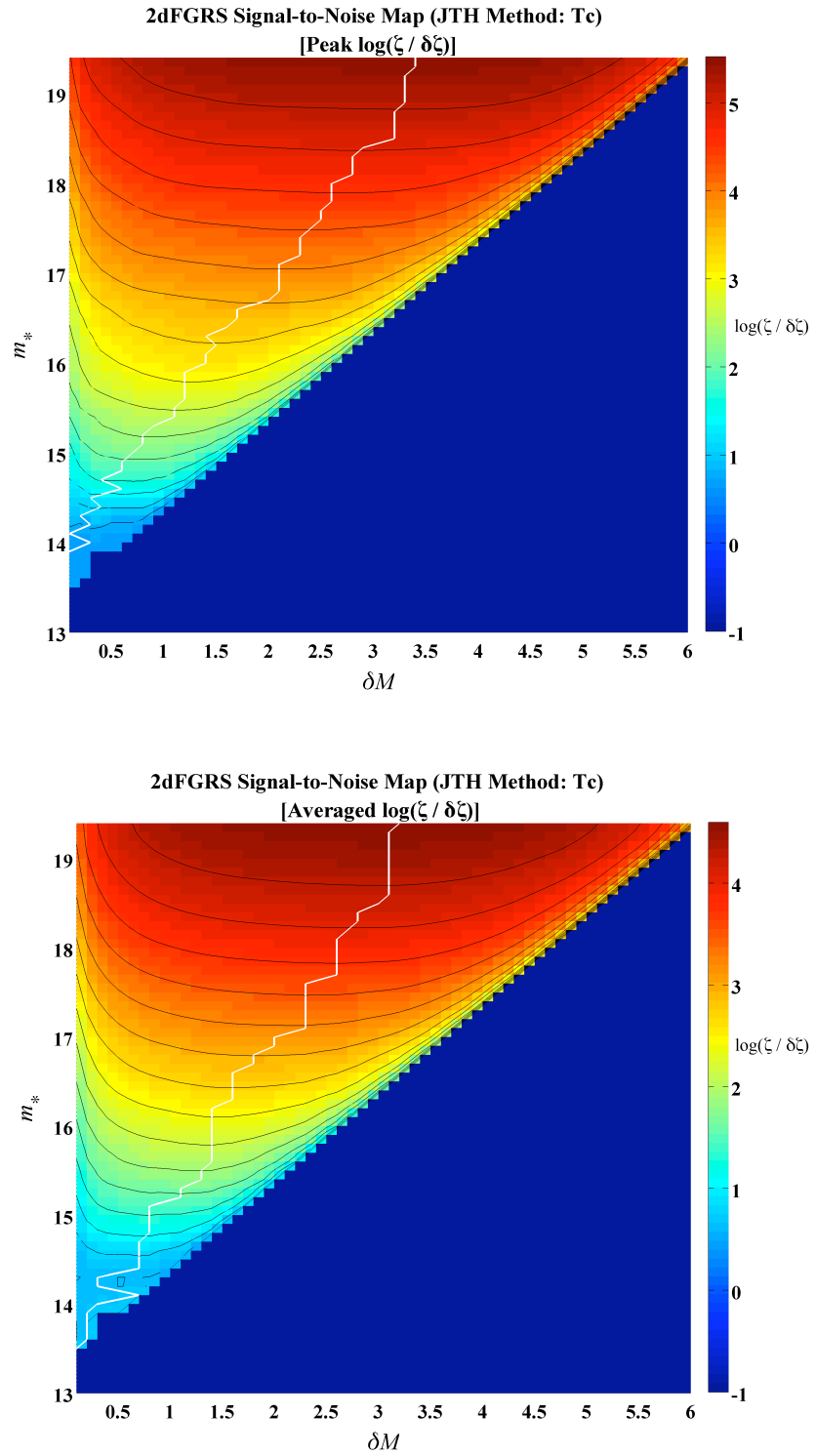


Figure 7.30: 2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.

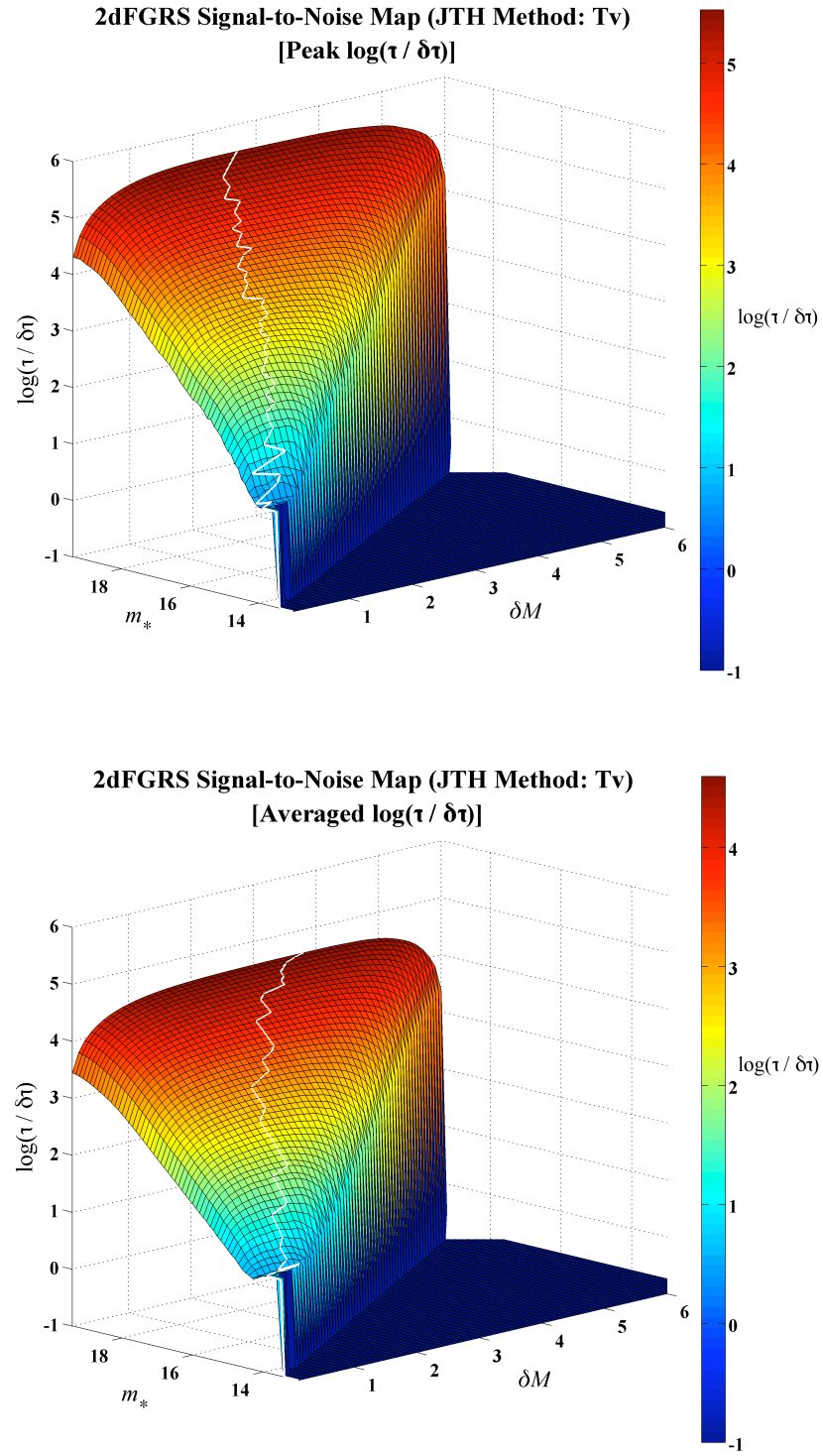


Figure 7.31: 3D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.

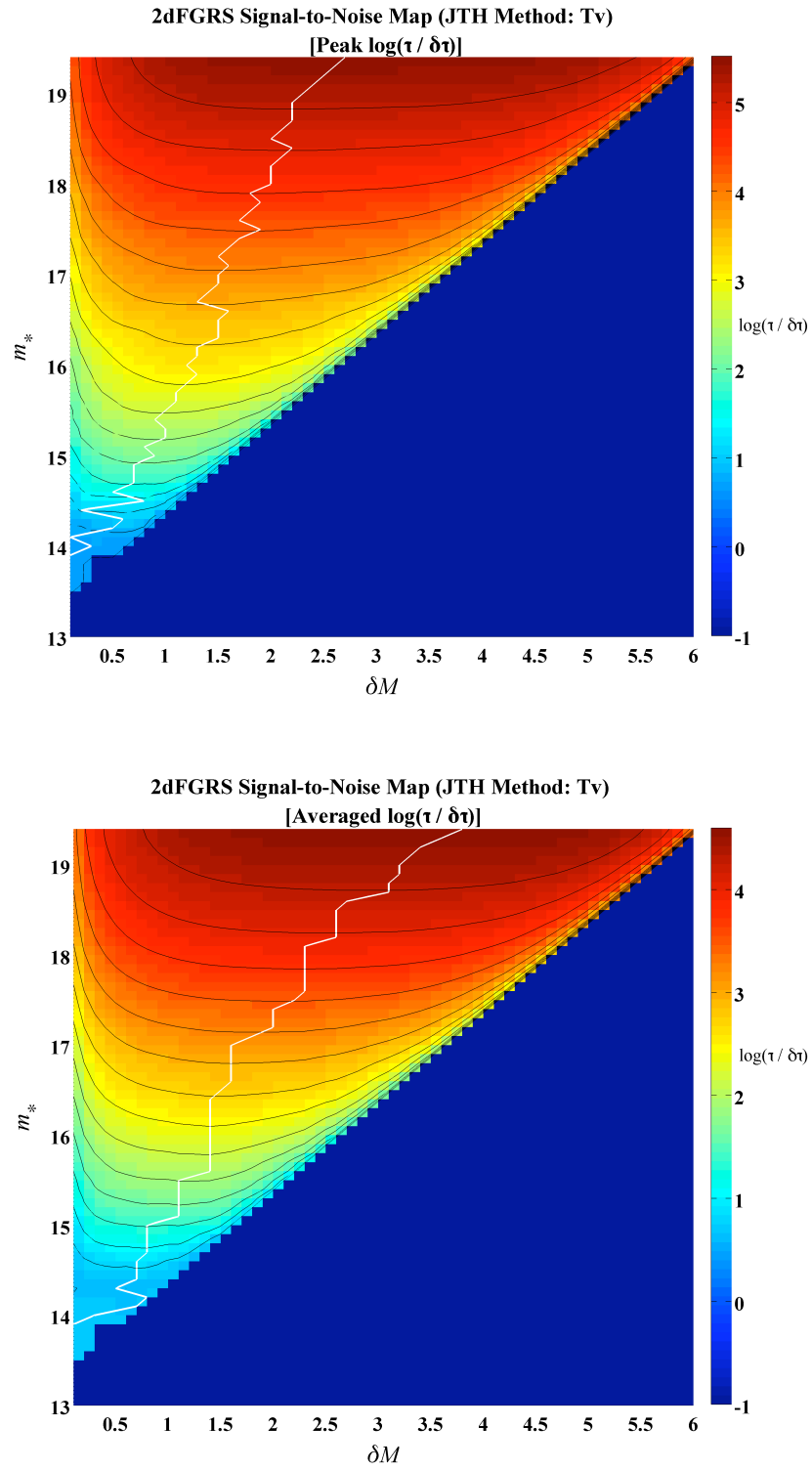


Figure 7.32: 2D representation of the SDSS signal-to-noise map derived from the application of the JTH T_v statistic.

7.3 Conclusions

This chapter has been largely concerned with specific details of the JTH method that have arisen from analysis of survey samples from 2dFGRS and CCQLG. In the first half of the chapter we explored two effects of the JTH method that, if not accounted for properly, could lead to inaccurate conclusions. The first of these concerned a resulting flat-lining of the T_c and T_v statistics beyond the survey apparent magnitude limits. We concluded that this effect manifested if the respective widths δZ and δM applied to the statistics were too small resulting in a shot-noise dominated scenario. Moreover, for each survey that we explored - MGC, 2dFGRS, SDSS and CCLQG - the resulting flat-line would occur at the exact limit of the survey i.e. the faintest galaxy in the data-set. However, as we have identified, particularly in the 2dFGRS and CCLQG surveys, this is not necessarily where the true completeness limit lies. Indeed, for these particular survey samples we have observed that they do not have well defined sharp apparent magnitude limits and therefore for a δZ and δM that is too small, there is no possibility of accurately determining the true limit of the survey.

Conversely, the second effect arose from the case where δZ and δM become sufficiently large that the respective T_c and T_v statistics both indicate a range where the true apparent magnitude limit may lie. In order to overcome these issues and improve and optimise our estimator, we are currently working on estimating ζ and τ based on their signal-to-noise ratio calculated from the respective areas, S_1, S_2 and S_3, S_4 . Thus far we have generated s/n maps for each of the surveys already examined by applying the JTH method for incremental range of δZ and δM . These maps have allowed us to trace the optimal s/n as a function of the trial apparent magnitude limit, m_* and the widths, δZ and δM . Our continuing work with this will lead to using these maps to optimise T_c and T_v based on a constant or minimum s/n level. Furthermore, we acknowledge that the implementation of either the original JTH or this improved optimised approach still leaves scope for manipulation to gain the best possible value for determining the true apparent magnitude limit. Our future work will therefore develop a complimentary extension to our method that will provide a comprehensive error estimate for both T_c and T_v .

Chapter 8

Creating Mock Galaxy Catalogues

“I have no data yet. It is a capital mistake to theorise before one has data. Insensibly one begins to twist facts to suit theories instead of theories to suit facts.”

Mr Sherlock Holmes - *A Scandal in Bohemia* (1892)

So far we have seen how real survey data has been instrumental in developing and implementing our statistical methods. However, real data carries the possibility that systematics, from sources such as errors in magnitude, can creep in and prompt us to draw wrong conclusions. For the work that follows in the next chapter, we needed to have ‘clean’ catalogues where we could mimic effects that appear in nature such as evolution without the risk that the signatures of these effects might be masked or distorted.

In this chapter we turn our attention to the production of Monte Carlo, or mock, simulated galaxy redshift surveys. Simulating galaxy surveys provides us with a controlled testing ground that, in our case, will facilitate the development of our current statistics for the study of evolution. We provide the details of the procedure adopted which allows us to efficiently mock surveys such as, MGC, 2dFGRS and SDSS (Early Types) before going on to compare these mocks to the actual survey data and finally perform a completeness analysis.

8.1 Different Methodologies

We firstly consider two alternative approaches to the generation of mock galaxy catalogues, in order to highlight some of the pertinent issues.

8.1.1 Theoretician's approach

The first is a standard approach to generate mock catalogues which could be thought of as the ‘theoreticians approach’. If one knows enough of galaxy formation and has access to state of the art computers to run simulations, then one can generalise this method as one of simulating,

$$\begin{aligned} p(M_{abs}, z) &= \int p(M_{abs}, z, m_{DM}) dm_{DM} \\ &= \int p(z|M_{abs}, m_{DM})p(M_{abs}|m_{DM})p(m_{DM}) dm_{DM} \end{aligned} \quad (8.1)$$

Where $p(M_{abs}|m_{DM})$ is called the *mass-to-light ratio*, m_{DM} denotes the distribution of dark matter halos, z the redshift, M_{abs} is absolute magnitude, and $Pr(z|m_{DM})$ expresses our knowledge regarding galaxy formation and evolution.

We can cast Equation 8.1 in the framework of the ‘Halo Occupation Distribution’ (HOD) (see [Seljak, 2000](#); [Berlind and Weinberg, 2002](#)). HOD describes the bias of a given class of galaxies by specifying the probability of finding N such galaxies in a dark matter halo,

$$\begin{aligned} p(N, z) &= \int p(N, z, m_{DM}) dm_{DM} \\ &= \int p(N|z, m_{DM})p(z|m_{DM})p(m_{DM}) dm_{DM} \end{aligned} \quad (8.2)$$

where, m_{DM} is the virial mass of a dark matter halo containing N galaxies (see above). From either of these starting points, we can describe mock generation in two main steps:

- Draw from a cosmological N-body simulation, generated in a Λ CDM framework which samples the mass distribution of dark matter halos to produce a catalogue that has realistic clustering. Galaxies are ‘assigned’ to each dark matter halo and then one chooses a location and orientation of the observer within the simulation that represents the local Universe.

- From a given redshift survey, one adopts a luminosity function, typically the Schechter function, that has parameters, M_* , α , and Φ_* , which are inferred from observations and should represent the present day luminosity function. Finally, loop over all the galaxies in the simulation and select the galaxies that are consistent with the survey selection function, at which point an apparent magnitude is generated to be consistent with the redshift of the galaxy from the simulation.

8.1.2 Observer’s approach

If the first method for generating mocks can be characterised as the theoretician’s approach, the latter method can conversely be thought of as the ‘observers approach’. It is this method that we have opted to use for the remainder of this thesis and it differs from the ‘theoretician’s approach’ in one crucial respect. Instead of using an N-body simulation, to generate a realistic spatial distribution of dark matter halos, which we then ‘populate’ with galaxies, we chose to use the observed redshift distribution from a known galaxy survey, which allows us to re-write Equation 8.1 in a simpler form,

$$\Pr(M_{abs}, z) = \Pr(M_{abs}|z)P(z) = \Pr(z|M_{abs})\Pr(M_{abs}) \quad (8.3)$$

The first of the equalities has the advantage of being able to replicate the observed redshift distribution *exactly*, whilst the second mimics the luminosity function exactly. In both cases they do not demand the scale of computing power required from N-body simulations. Conversely, we could use the observed magnitude distribution from a known survey. However, whichever observed distribution we choose, this method obviously lacks the scope for producing multiple realisations of that distribution where the effects of cosmic variance (particularly in smaller volumes at high redshift) for one realisation make it difficult to draw statistical conclusions (Somerville et al., 2004). It becomes clear that ultimately the way to overcome this limitation is to adopt the ‘theoretician’s approach’ above where one can effectively generate as many ‘Universes’ as one requires. However, for the purposes of our analysis the observer’s approach will be more than sufficient.

8.2 Which Frame?

Having chosen to generate mocks from an observed redshift distribution we now consider the way in which we will sample the absolute magnitudes. This leads us to consider two possible ‘reference frames’ in which our magnitudes may be selected. The first we call the ‘observer’s frame’ and the other we have termed the ‘galaxy’s frame’.

The Observer's Frame: represents a realisation of the actual observed apparent (or absolute) magnitude distribution of galaxies from a survey prior to any corrections (such as evolution and/or k-correction). In practice, the process in which one obtains a real distribution in this frame can be summarised as follows.

1. **Galaxy detection:** An observational cosmologist will measure apparent magnitudes of galaxies in a portion of the sky out to a faint limiting apparent magnitude, $m_{\text{lim}}^{\text{f}}$ imposed by the physical limitations of the telescope. Limitations on the CCD instrumentation, where effects from pixel saturation due to very bright objects also imposes a bright limit to the survey, $m_{\text{lim}}^{\text{b}}$. The data reduction side of this process is by no means a trivial task and there are several effects that hinder the process of obtaining a magnitude complete catalogue. For example, objects that are close to bright stars may be missed as well as objects that either lie at the edge of the image or at a defected part of the detector. Moreover, galaxies with the same surface brightness may or may not be detected depending on their shape and overall extent: a compact object is more likely to have enough pixels above the detection limit than a very diffuse galaxy of the same brightness.
2. **Extinction correction:** Once the observing run is complete, it is then necessary to correct the magnitudes for galactic extinction using extinction maps as in for example, [Schlegel et al. \(1998\)](#). Although there is of course also intrinsic galactic extinction for each targeted galaxy, it is a very difficult quantity to measure and therefore generally ignored. Other related effects that must be carefully accounted for include atmospheric dust and varying sky brightness.
3. **Measure redshifts:** At this stage the catalogue consists of measured magnitudes and sky positions only, without their 3D spatial redshift distribution. Therefore, each galaxy is then targeted to obtain a measure of its redshift either spectroscopically or photometrically. Although the majority of surveys obtain redshifts spectroscopically, as it is far more accurate than the photometric counterpart, the main drawback of spectroscopic redshifts arises from spatial limitations. For example, in a region that has a high density of objects, a lot of galaxies may be missed since there is a physical spatial limit on how close the fibres can be placed close enough together.
4. **Impose a magnitude cut:** It is often necessary to impose a magnitude limit that is brighter than the detection limit of the instrument. This is commonly due to effects such as the required signal-to-noise of the spectra. If one is obtaining

redshifts spectroscopically then effects caused by the number of fibres, sensitivity and field of view of the spectrograph will also play a key role in the magnitude cut of the survey. Similarly, for surveys that use photometric redshifts, limitations are imposed by the accuracy of the photometry.

5. **Combine the data:** Finally, both the imaging and the redshift data are combined to form the final catalogue. What we would have now represents the publicly released data that, for the surveys we have already analysed, has an M - Z distribution that resembles e.g. Figures 3.4 (MGC), 4.2 (2dFGRS) and 4.20 (SDSS-Early Types).

At this stage the magnitudes have not been corrected to take into account effects such as k -correction and possible evolution within the survey. The data therefore represents the observer's distribution, hence the term, 'observers frame'.

The Galaxy Frame: We want to compare galaxies in the rest (or the galaxy's) frame, or equivalently at $z = 0$. Since we generally observe each galaxy through a single bandpass, we only see a fraction of its total spectrum which is redshifted into the observer's frame. Therefore, each galaxy's magnitude requires to be k -corrected to take this effect into account. Furthermore, the further away a galaxy is from us, the more its total magnitude will be affected from evolutionary effects. If all these corrections have been accurately accounted for, in principle one should be able to determine a present day universal luminosity function from a given survey.

Therefore, to summarise in the context of the variables we have been using throughout this thesis,

$$\text{Galaxy's Frame} \begin{cases} m_{\text{corr}} = m_{\text{obs}} - k(z) - e(z) - A_v(l, b), \\ M_{\text{corr}} = m_{\text{corr}} - 5 \log_{10}(d_L) - 25, \\ Z = 5 \log_{10}(d_L) + 25. \end{cases}$$

$$\text{Observer's Frame} \begin{cases} m = m_{\text{obs}}, \\ M = m - 5 \log_{10}(d_L) - 25, \\ Z = 5 \log_{10}(d_L) + 25, \end{cases}$$

where $k(z)$ is the redshift dependent k -correction, $e(z)$ is the redshift dependent evolutionary model correction, m_{obs} is the apparent magnitude of a galaxy in the observers

frame prior to the application of these corrections and m_{corr} is the final corrected magnitude that brings the galaxy from the observer's frame to the galaxy (or rest) frame.

8.3 Our Mock Recipe

Since we chose to adopt the observer's approach for generating our mock catalogues, it was a natural step to implement the redshift distributions from the MGC, 2dFGRS and SDSS (Early-Types) samples that we have already analysed. Moreover, this allowed us the opportunity to adopt similar, and in the case of MGC, the same LF parameters already derived by the respective teams. The process in which we create a mock catalogue based on the Universal LF derived from the corrected magnitudes for each survey is now described.

8.3.1 Sampling the magnitudes

The LF from which we are sampling the magnitudes for the MGC and 2dFGRS catalogues is the widely used Schechter function of [Schechter \(1976\)](#). For the SDSS (Early Types) we sample from a Gaussian function as applied in [Bernardi \(2003b\)](#).

Initially, it was our intention for the MGC and 2dFGRS mocks to sample luminosities, L , directly from from a Schechter function and then convert to them to absolute magnitudes, M . Such an approach would therefore involve generating a Monte Carlo sample drawn from a gamma distribution since from Equation 1.18 on page 23 we have:

$$\Gamma(\alpha) = \int_0^{\infty} t^{\alpha-1} e^{-t} dt, \quad (8.4)$$

$$\text{where, } t = \left(\frac{L}{L_*} \right),$$

then convert them to absolute magnitudes,

$$M = M_* - 2.5 \log \left(\frac{L}{L_*} \right). \quad (8.5)$$

Sampling magnitudes in this way, directly from the Schechter function would both be computationally efficient and elegant. However, the published value of α in surveys such as the 2dFGRS is negative with a typical value of $\alpha = -1.21$ as in [Norberg \(2002a\)](#). After searching extensively through the literature, although there are gamma function

random number generators, it would appear there are none that sample directly from a gamma distribution for negative values of α .

Having to reluctantly abandon this direct sampling approach we used the more conventional method of the probability integral transform, or as it is also known, ‘inverse transform sampling’. The transformation $F(X)$, of a random variable, X , with an invertible cumulative distribution function, F , to a uniform distribution between $[0,1]$ allows us to generate a random sample drawn from the PDF of X . In the present case our random variable is absolute magnitude, M , and our function is the luminosity function, $\Phi(M)$. This requires us to firstly convert the LF from luminosities, L to absolute magnitudes, M . From Equation 1.18 we have:

$$\Phi(M)dM = \Phi(L)dL = \phi_* \left(\frac{L}{L_*} \right)^{-\alpha} \exp \left(-\frac{L}{L_*} \right) dL \quad (8.6)$$

The relation between absolute values of luminosities and magnitudes is given by,

$$\frac{L}{L_*} = 10^{-0.4(M-M_*)} = \exp[-0.4 \ln 10(M - M_*)] \quad (8.7)$$

It follows that,

$$\frac{dL}{dM} = -0.4L_* \ln 10 \{ \exp[-0.4 \ln 10(M - M_*)] \} \quad (8.8)$$

Substituting Equations 8.7 and 8.8 into Equation 8.6 gives:

$$\Phi(M) = \phi_* \exp[-c\alpha(M - M_*)] \exp[-\exp(-c(M - M_*))] \exp[-c(M - M_*)] \quad (8.9)$$

$$= \phi_* \exp \{ [-c(\alpha + 1)(M - M_*)] - \exp[-c(M_* - M)] \} \quad (8.10)$$

where, $c = 0.4 \ln 10$.

Sampling from a Gaussian function: To construct a mock catalogue for the SDSS (Early Types) we have sampled magnitudes from a Gaussian function of the form given by Bernardi (2003b),

$$\Phi(M_i, z_i | M_*, \sigma_M) = \frac{\phi_*}{\sqrt{2\pi\sigma_M^2}} \times \exp \left(-\frac{[M_i - M_*]^2}{2\sigma_M^2} \right), \quad (8.11)$$

where M_* is once again the characterisitic absolute magnitude, ϕ_* is the comoving number density of the galaxies, and σ_M is the variance.

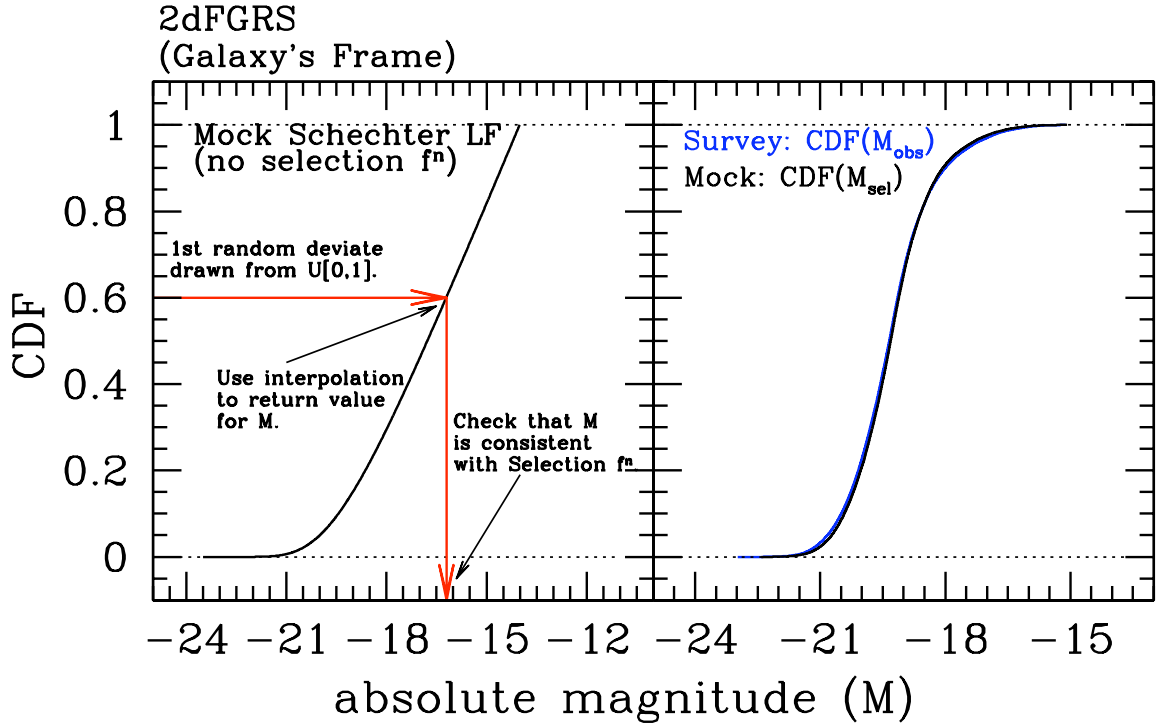


Figure 8.1: Schematic example of how we sample the absolute magnitudes for our mock catalogue from the CDF of a Schechter function . We have used the 2dFGRS for this example. The left-hand panel shows M_{CDF} vs the normalised CDF of the Schechter function. Using a random number generator, we produce a random deviate from a uniform random distribution between $[0,1]$. We then used cubic spline interpolation to return a value for M_{samp} and check to see if its value lies within both absolute and apparent magnitude limits. If it is, we accept the value, if it is not, we reject it and resample. This process is repeated until we have a full catalogue based on the number of redshifts from the 2dFGRS. We can then compare the CDF's from our new mock catalogue with the actual observed absolute magnitude distribution as shown on the right hand panel. It is clear that both CDF's match very well indicating that our selection procedure is working.

8.3.2 Final selection

The selection procedure for the magnitudes can be broken down into the following steps which detail how we can convert from the galaxy's frame (G-Frame) to the observer's frame (O-Frame).

1. We firstly create a CDF from the appropriate LF by generating a series of absolute magnitudes, M_{CDF} , in a predetermined range $M_{\text{max}} > M_{\text{CDF}} > M_{\text{min}}$ that encompasses the range of the observed survey we are attempting to mock. The LF parameters used are those already derived from the corrected magnitudes of the real survey in the G-Frame. We then normalise the CDF as illustrated in the

left hand panel in Figure 8.1.

2. For each redshift, z_i in our catalogue, we sample absolute magnitudes from the CDF by generating a random uniform deviate between $[0, 1]$ and then use cubic spline interpolation to determine a corresponding sampled value M_{samp} (see Figure 8.1). We then convert this to a sampled apparent magnitude $m_{\text{samp}}(z_i)$ using,

$$m_{\text{samp}}(z_i) = M_{\text{samp}}(z_i) + Z_i,$$

where the distance modulus, Z , is calculated directly from our observed redshift distribution assuming a cosmological model.

3. We have now sampled magnitudes from a LF that represents the corrected magnitudes in the galaxy's frame. At this stage we must impose our selection criteria which ensures the value, $M_{\text{samp}}(z_i)$, satisfies the conditions where by the galaxy lies within both the apparent *and* absolute magnitude limits. However, these limits are initially defined by the raw uncorrected magnitudes in the observer's frame. We therefore, need to move from the observer's frame to selecting galaxies in the galaxy's frame which is achieved satisfying the following conditions,

$$\text{Condition 1} \begin{cases} m_{\text{lim}}^{\text{f}}(\text{mock}) &= m_{\text{lim}}^{\text{f}}(\text{survey}) - k(z) - e(z) \\ m_{\text{lim}}^{\text{b}}(\text{mock}) &= m_{\text{lim}}^{\text{b}}(\text{survey}) - k(z) - e(z) \end{cases} \quad (8.12)$$

$$\text{Condition 2} \begin{cases} M_{\text{lim}}^{\text{f}}(\text{mock}) &= m_{\text{lim}}^{\text{f}}(\text{mock}) - 5 \log_{10}[d_L(z_i)] - 25 \\ M_{\text{lim}}^{\text{b}}(\text{mock}) &= m_{\text{lim}}^{\text{b}}(\text{mock}) - 5 \log_{10}[d_L(z_i)] - 25 \end{cases} \quad (8.13)$$

If the galaxy satisfies both conditions it is accepted and Steps 2 and 3 are repeated for the next galaxy at redshift z_i . Otherwise it will be rejected and the absolute magnitude will be resampled as per Step 2 until conditions 1 and 2 are satisfied. By including both evolutionary and k-correction into our selection criteria we are creating a mock catalogue in the **galaxy's** frame (G-Frame). That is, the magnitudes selected in our mock represent the corrected magnitudes of the survey we are simulating. If, on the other hand, we wish to mock the uncorrected survey data which represents the **observer's** frame (O-Frame) we have to introduce a final step.

4. To convert the final selected G-Frame mock magnitudes, $m_{\text{sel}}^{\text{GF}}$ and $M_{\text{sel}}^{\text{GF}}$, to the

observer's frame we simply apply the k - and evolution-corrections such that,

$$m_{\text{sel}}^{\text{OF}} = m_{\text{sel}}^{\text{GF}} + k(z) + e(z) \quad (8.14)$$

$$M_{\text{sel}}^{\text{OF}} = M_{\text{sel}}^{\text{GF}} + k(z) + e(z), \quad (8.15)$$

where $m_{\text{sel}}^{\text{GF}}$ and $M_{\text{sel}}^{\text{GF}}$ are $m_{\text{samp}}(z_i)$ and $M_{\text{samp}}(z_i)$ respectively provided that the sampled galaxy has passed the selection criteria at step 3.

We now have a mock catalogue which is a true representation of the actual survey data (before corrections). In our example in Figure 8.1 we have mocked a sub selection of the 2dFGRS. The right hand panel compares the CDF's of both the actual observed absolute magnitudes for the survey to our mock magnitudes. As expected, these show close agreement.

Depending on the frame in which we are wishing to present our mock catalogue, we finally loop over all galaxies in our redshift catalogue until all the magnitudes have been selected.

8.4 Mocking MGC, 2dFGRS and SDSS (Early Types)

Now that we have established our procedure for generating mock surveys, we apply it to the three surveys that have already played a key role in our understanding and extension of the T_c and T_v statistics - MGC, 2dFGRS and SDSS (Early Types) - and compare the distribution of the mock surveys that we generate to the actual survey data.

8.4.1 MGC

For the MGC redshift distribution we have used the exact selection criteria we applied in Chapter 3, section 3.2.2 on page 60, giving us a total of 7878 galaxies in the range $0.013 < z < 0.18$. We recall that for this survey a global pure luminosity evolution correction was applied given by,

$$E(z) = -\beta \times 2.5 \log_{10}(1 + z_i) \quad (8.16)$$

where a value for the evolution parameter, $\beta = 0.75$ was found to be suitable. The k -corrections as applied in Driver et al. (2005) (D05) were derived on an individual galaxy basis utilising a fitting procedure to 27 spectral templates from Poggianti (1997),

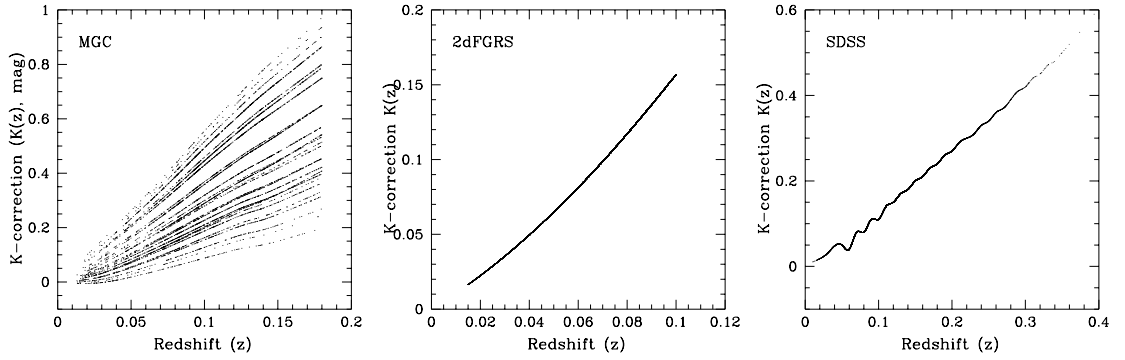


Figure 8.2: Plots tracing the k -correction models as applied to MGC, 2dFGRS and SDSS (Early Types). The left-hand plot traces the 27 k -correction spectral templates as applied to the MGC survey data by [Driver et al. \(2005\)](#). The middle panel shows the global k -correction for our 2dFGRS which was derived by [Norberg \(2002a\)](#), and the right panel is k -corrections applied to our SDSS survey sample as derived in [Bernardi \(2003a\)](#).

as illustrated in the left-hand plot in Figure 8.2. If our sole goal was to produce a completely accurate representation of the survey data, it would obviously be necessary also to simulate the varied galaxy type and apply the same k -correction procedure as in D05 for each realisation. However, for the purposes of our analysis in next chapter it is sufficient to apply the same k -correction values as derived for the observed survey data.

In D05 the authors recover Schechter LF parameters from a joint luminosity-surface-brightness step-wise maximum-likelihood procedure where the above selection effects have been taken into account. Their fitted parameters are: $\phi^* = (0.0177 \pm 0.0015)h^3\text{Mpc}^{-3}$, $M_{B_{MGC}}^* - 5\log h = (-19.60 \pm 0.04)\text{ mag}$ and $\alpha = -1.13 \pm 0.02$. Using the notation from Equations 8.12 and 8.13, we have set the faint and bright apparent magnitude limits to $m_{\text{lim}}^f(\text{survey}) = 20.0\text{ mag}$ and $m_{\text{lim}}^b(\text{survey}) = 13.64\text{ mag}$. From this, we have generated 1000 MGC realisations in both the O-Frame and the G-Frame. Figure 8.3 shows the apparent and absolute magnitude distributions for the real survey compared to our mocks. The blue shaded regions represent the actual survey whilst the red are the mock distributions. In the top panels we can see the resulting apparent magnitude distribution where the left-hand panel represents the G-Frame and the right-hand panel is the O-Frame. Comparing our mocks to the survey we can see that they agree very well with only slight discrepancies towards the peaks of the distribution, a result mirrored in the corresponding absolute magnitude distributions shown in the bottom panels.

In the G-Frame panel for the apparent magnitudes we observe a systematic change

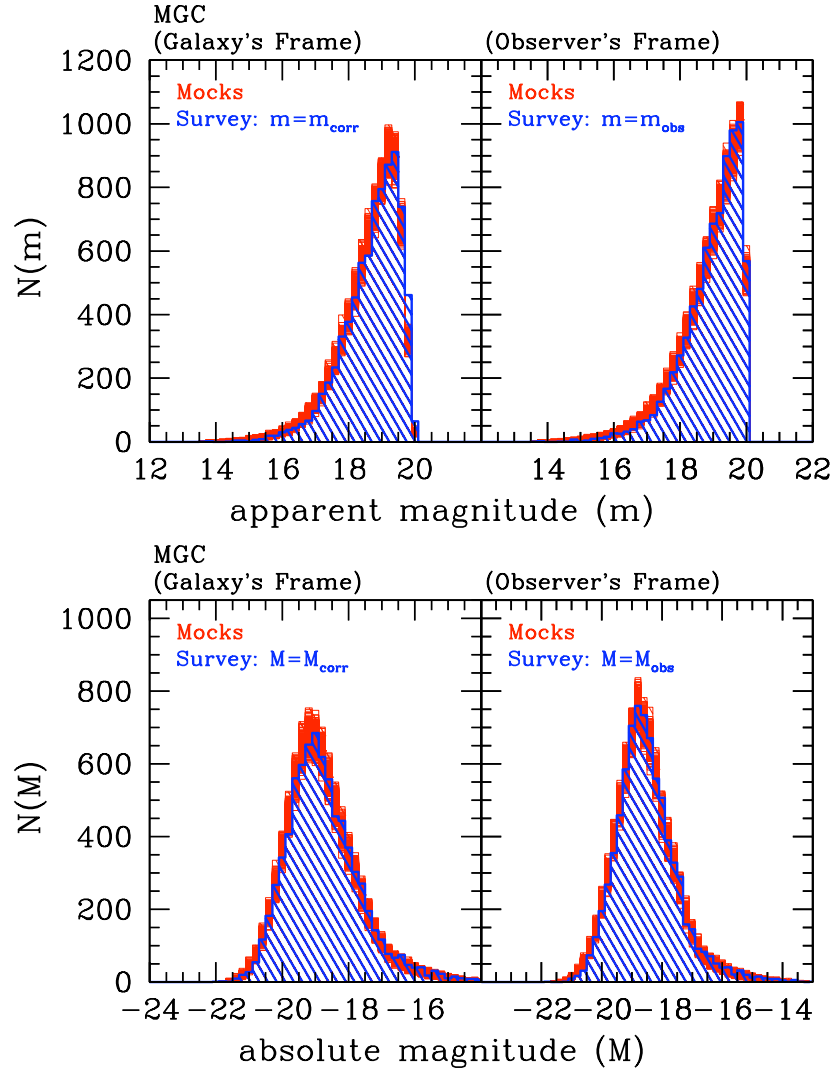


Figure 8.3: Apparent and absolute magnitude distributions comparing our MGC mocks to the real data. The top panels show the histograms of apparent magnitudes for the actual MGC survey (blue region) and 1000 of our mock catalogues (red regions). The left-hand of the two panels considers this distribution in the galaxy frame where, in the case of the survey data, the magnitudes have been corrected for evolution, k -correction and galactic extinction. Consequently, the right-hand panel shows the observer's frame that represents the distribution of magnitudes before any of the aforementioned corrections have been applied to data. The bottom panels mirror the top for the distribution of absolute magnitudes. All four plots demonstrate that our mocks show close agreement with the actual MGC data showing only slight discrepancies towards the peak of the distributions. This is most likely due to the fact that our mocks do not account for any galactic extinction modelling in our selection function.

in the faint magnitude limit compared to the O-Frame for the uncorrected data. This is caused by the $(k+e)$ -corrections that have been added to the magnitudes which

causes the apparent magnitude limit to curve according to these models, this effect is shown in M - Z distributions in Figure 8.6-top . In this figure we have selected a mock at random (left-hand panel) to compare it to the survey (right-hand panel). Since no corrections have been added to the distance modulus distribution, we observe the apparent magnitude limit to curve in such a way modelled by the $(k + e)$ -corrections, shown by the black points in diagram, relative to the uncorrected magnitudes shown by the red points. If we recall our initial completeness analysis of the actual MGC and 2dFGRS survey's, we note that this behaviour was not observed, since we followed the procedures laid out in Rauzy (2001), where, in addition to correcting both the apparent and absolute magnitudes, the distance modulus was also corrected. This approach ensures that the apparent magnitude limit (faint or bright) will always remain straight. We explore both of these approaches and its effects on our statistical analysis in greater detail in Appendix-A.

8.4.2 2dFGRS mocks

For the 2dFGRS mocks we decided to use a smaller subsample than was used in our initial completeness analysis in Chapter 4 in order to maximise the computational time. We have therefore selected galaxies between $0.015 < z < 0.1$ and within the apparent magnitude range $12.0 < m < 19.0$. This now gives a more manageable subsample of 66696 galaxies. In Norberg (2002a) they derived the following Schechter LF parameters: $\phi^* = (0.0161 \pm 0.08)h^3\text{Mpc}^{-3}$, $M_* - 5 \log h = (-19.66 \pm 0.07)$ mag and $\alpha = -1.21 \pm 0.02$ from a sample of 110,500 galaxies out to $z < 0.2$ between $17.0 < m < 19.0$. We therefore do not expect these LF parameters to completely accurately describe our subsample, however, we can use them as a guide to draw a suitable mock magnitudes. We also use the global $(k+e)$ -correction (see middle panel in Figure 8.2) derived in Norberg (2002a) and given by,

$$k(z) + e(z) = (z + 6z^2)/(1 + 20z^3) \quad (8.17)$$

Figure 8.4 shows our initial magnitude distributions for a trial 2dFGRS mock which samples magnitudes using the exact LF parameters as stated above. The apparent magnitudes are shown on the left-hand panel and the absolute magnitudes are on the right. Unsurprisingly, when we compare our mock (shown in red) with the actual survey sample we are using (shown in blue) they do not match. However, by changing the value of M_* from -19.66 mag to -20.10 mag we now recover magnitude distributions shown in Figure 8.5 which are much closer to the survey distribution. Once again, in Figure 8.6-middle, we can view the M - Z distributions for one of the mocks (left) compared to

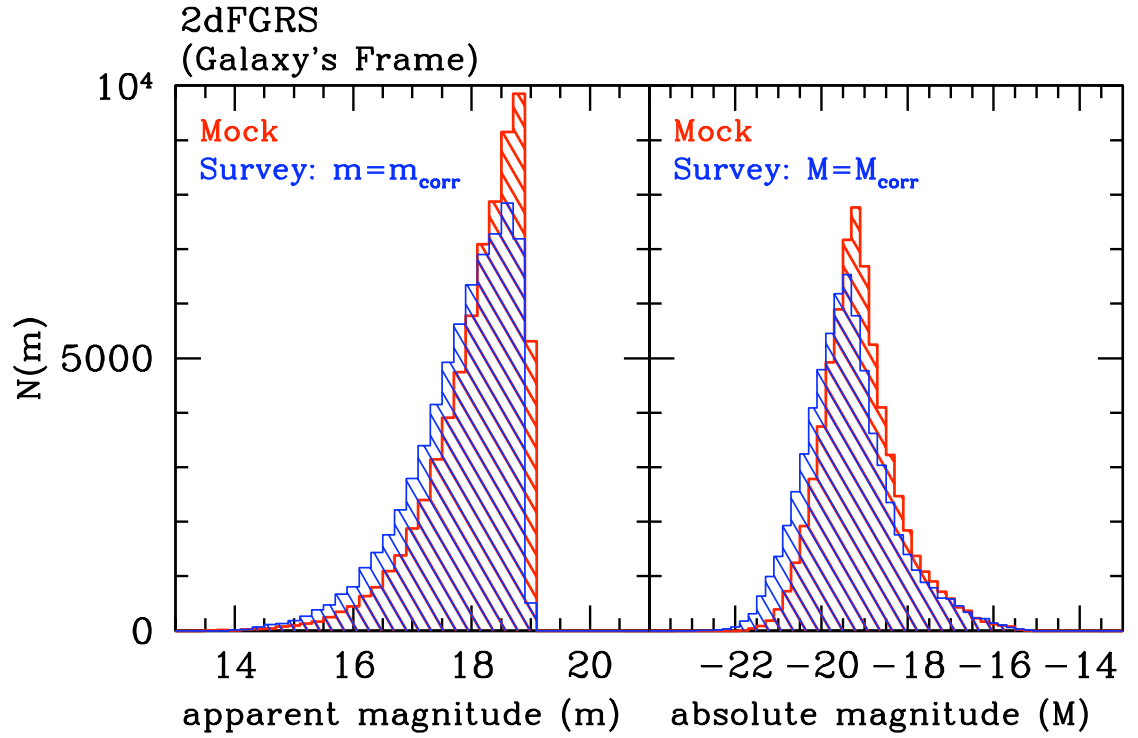


Figure 8.4: Apparent and absolute magnitude distribution for a trial 2dFGRS mock compared to the observed survey subsample we have selected. The red distribution in both panels shows our first trial mock adopting the exact LF parameters as in [Norberg \(2002a\)](#). As expected, since our subsample does not match the same survey sample in the paper, both the absolute and apparent magnitude distributions do not match the observed magnitudes of the survey shown as the blue shaded regions.

the survey (right). In both cases they illustrate the curved apparent magnitude limit (G-Frame) for the corrected magnitudes and the straight magnitude limit (O-Frame) for the uncorrected magnitudes.

8.4.3 SDSS Mocks

Creating suitable SDSS mocks based on the Early Types catalogue uncovered the same issues with sampling the magnitudes correctly as we saw with the 2dFGRS mocks. The LF adopted here was a Gaussian from [Bernardi \(2003b\)](#) where the parameters recovered by the SDSS team were based on an early sample of approximately 9000 galaxies in the redshift range $0.01 \leq z \leq 0.3$ given as $\phi^* = (0.0058 \pm 0.003)h^3\text{Mpc}^{-3}$, $M_* = -21.15$ mag and $\sigma_M = 0.84$. Since the redshift catalogue we are sampling from consists of 35421 galaxies out to $z < 0.4$ we realised that these published LF parameters were useful only as guide. After a little tweaking we found that $M_* = -20.80$ mag and

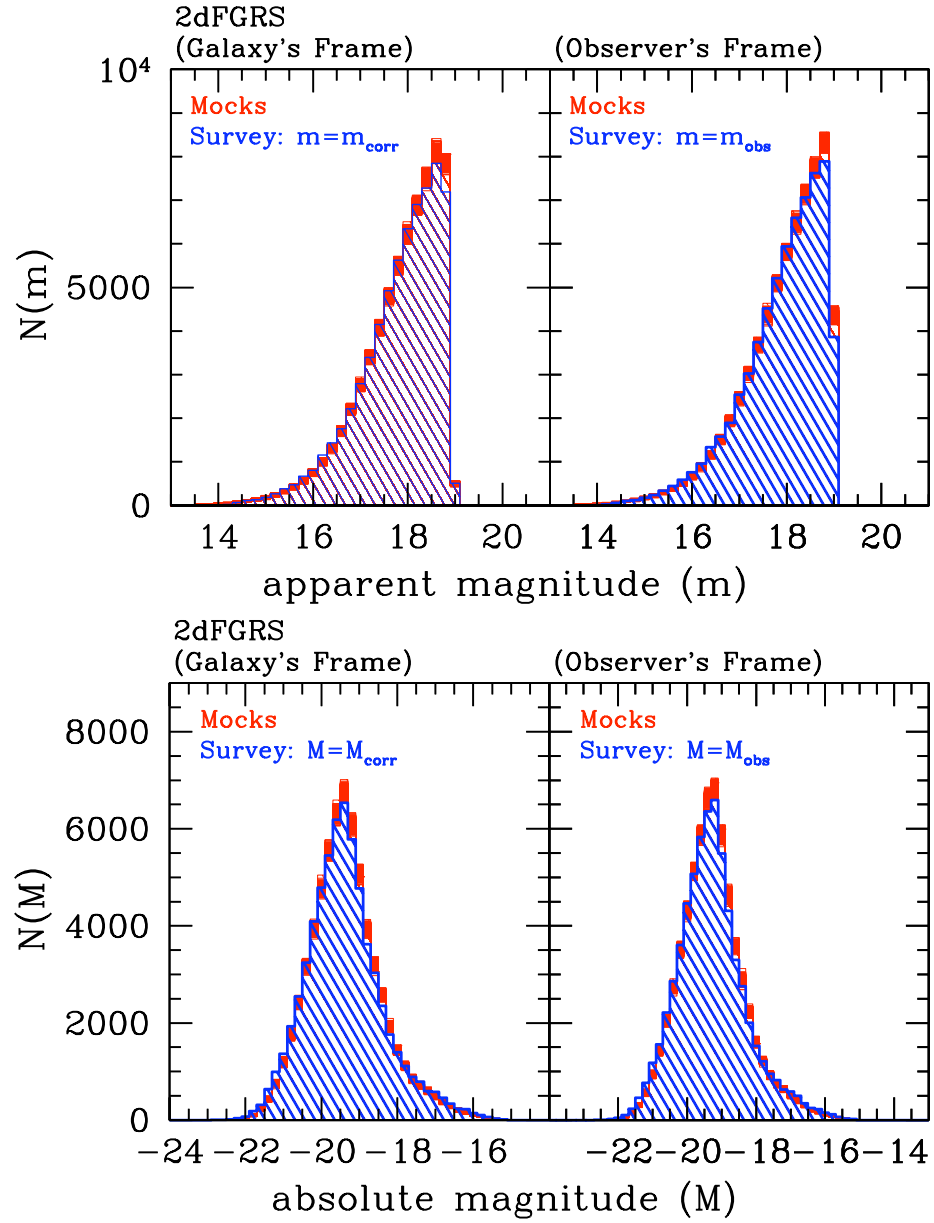


Figure 8.5: Histograms of apparent and absolute magnitudes comparing our 2dFGRS mocks to the real data. The top panels show the histograms of apparent magnitudes for the actual MGC survey (blue region) and 1000 of our mock catalogues (red regions). Although our LF parameters were based on those derived in [Norberg \(2002a\)](#), we found it necessary to vary M_* and α slightly since our subsample does not represent the same one used in that paper. All four plots demonstrate that our mocks show close agreement with the actual 2dFGRS data showing only slight discrepancies towards the peak of the distributions.

$\sigma_M = 0.87$ provided a reasonable fit for mock magnitude distributions as shown in [Figure 8.7](#).

We apply the same k -correction model as described in detail in Appendix-A of [Bernardi \(2003a\)](#). Figure 8.2-right traces the functional form of the template used. The bottom panels Figure 8.6 show the resulting M - Z distribution for one of the thousand mocks we generated compared to the survey sample data.

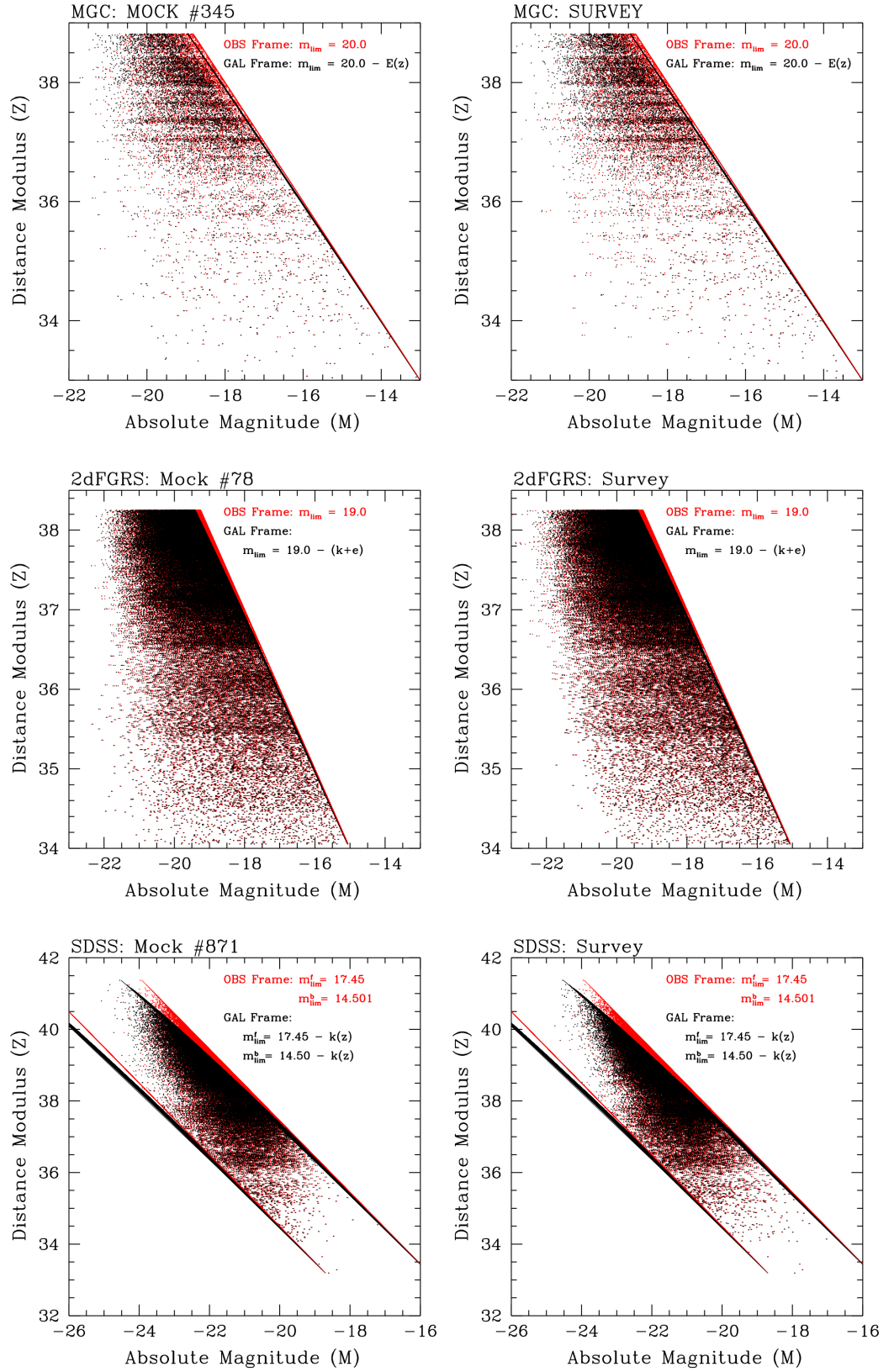


Figure 8.6: M - Z distributions for one of our MGC, 2dFGRS and SDSS mocks compared to their respective survey samples. The black points in all the panels show the corrected magnitudes, whilst the red points are the uncorrected magnitudes. What is evident with all the surveys is the apparent magnitude limit resulting from a modelled k - and e -corrections.

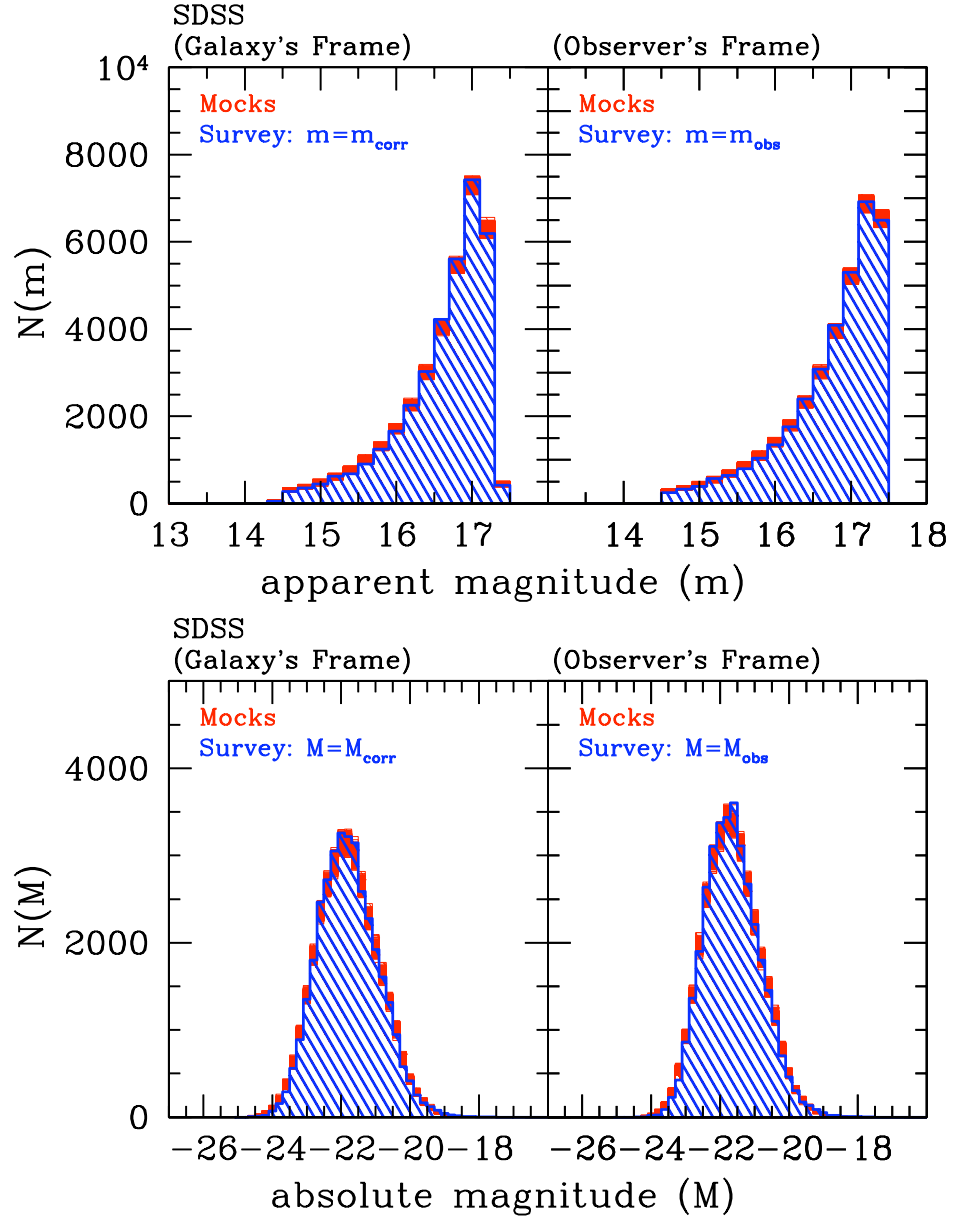


Figure 8.7: Apparent and absolute magnitude distribution for a trial SDSS (Early Types) mock compared to the survey subsample. The top panels show the histograms of apparent magnitudes for the actual SDSS survey (blue region) and 1000 of our mock catalogues (red regions). Similarly with 2dFGRS our LF parameters were based on those derived in [Bernardi \(2003b\)](#). However, the parameters derived were for an earlier data release with a smaller sample of only 9000 galaxies. To generate mocks that resembled the survey data we were considering, an $M_* = 20.8$ and $\sigma_M = 0.87$ was found to be adequate.

8.5 Completeness Analysis of the Mocks

In this section we turn our attention to the completeness of the mocks for all three surveys in both the galaxy's frame and the observer's frame. For MGC and SDSS we have run our completeness statistic for 1000 realisations. For the 2dFGRS mocks we tested total of 100 realisations due to the length of time which it takes to run the programs.

It is important to reiterate that when we apply T_c and T_v based the R01 and JTH methods in this chapter, the distance modulus, Z values remain uncorrected in the G-Frame with only the magnitudes being corrected.

8.5.1 Completeness of the MGC mocks

For the MGC mocks, we apply the R01 completeness procedure where we do not consider a bright limiting apparent magnitude. If we look at Figure 8.8 we observe what is at first glance a strange behaviour in T_c and T_v in the O-Frame (left-hand panel). The black and blue superimposed lines represent the respective 1000 mocks for T_c and T_v . We can clearly see a systematic upward trend in both statistics beginning at $m_* \sim 17.0$ mag and continuing up to the survey limit where T_v peaks at $\sigma \sim 7.0$ before dropping sharply as both statistics pass the survey limit at $m_{\text{lim}} = 20.0$. This kind of behaviour is perhaps not unexpected since the frame in which we are considering the magnitude distributions are uncorrected for effects such as evolution and k -correction. Therefore, we should expect the variables M and Z to be un-separable and consequently the assumption of separability to break down.

However, there does seem to be a region between the 3σ limits where the magnitudes have been sampled in such a way that render T_c and T_v to show consistent completeness up to the magnitude limit. This is in fact where we see the actual survey, shown in red, fluctuate. This perhaps an indication that effects from evolution and k -correction have a minimal impact on this survey due to the relatively shallow redshift range we are considering, but there exists certain sampled magnitude distributions where these effects can manifest.

For the G-Frame in the right-hand panel of Figure 8.8 we are now considering the corrected magnitude data. This time we T_c and T_v behave as one would expect for a separable, complete sample. Since, as we have already discussed, the Z -distribution remains fixed, this introduces an apparent magnitude limit that is curved according to the evolution and k -correction models applied to the data. In this case the overall trend is for the limit to become systematically brighter with increasing redshift. This

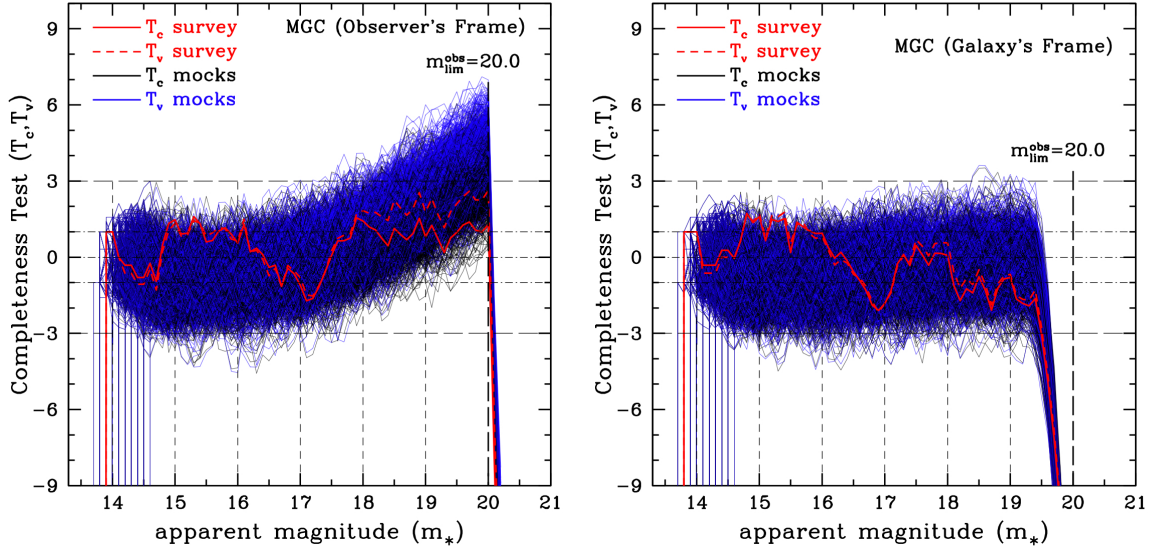


Figure 8.8: T_c and T_v results for MGC mocks in both the G-Frame (right-hand panel) and the O-Frame (left hand panel). In both cases we have applied the R01 method which does not account for a bright limit. The blue and black lines are the respective T_v and T_c results for a total of 1000 superimposed realisations, whereas the red solid and dashed lines are the respective T_c and T_v results for the real survey data. In the O-Frame, which represents the case of uncorrected magnitudes, we observe a systematic rise in both T_c and T_v for the mocks beginning at $m_* \sim 17.0$ mag continuing up to the survey limit where T_v peaks at $\sigma \sim 7.0$ before dropping sharply as both statistics pass the survey limit at $m_{\text{lim}}^{\text{obs}} = 20.0$. For the corrected magnitudes in the G-Frame on the right-panel we observe both statistics behaving as one would expect for a complete data-set. The mocks in this case indicate a range for the true magnitude limit between $19.3 \lesssim m_{\text{lim}}^* \lesssim 19.7$ with the real survey data indicating $m_{\text{lim}}^* \sim 19.5$ mag.

is reflected in the both T_c and T_v shown as a systematic drop below -3σ in the range $19.3 \lesssim m_{\text{lim}}^* \lesssim 19.7$. This result is confirmed by the actual survey data shown as red lines.

8.5.1.1 The distribution of T_c and T_v

Creating multiple realisations gives us an opportunity to examine some of the more fundamental aspects of our estimators. One such aspect is the resulting sampling distribution of T_c and T_v for any given m_* slice.

By their construction, both T_c and T_v should have a Gaussian sampling distribution with mean zero and variance equal to unity. This, therefore implies that if we have many realisations of a survey that is complete in apparent magnitude and apply the R01 method, we should expect for any given trial apparent magnitude limit, m_* , this

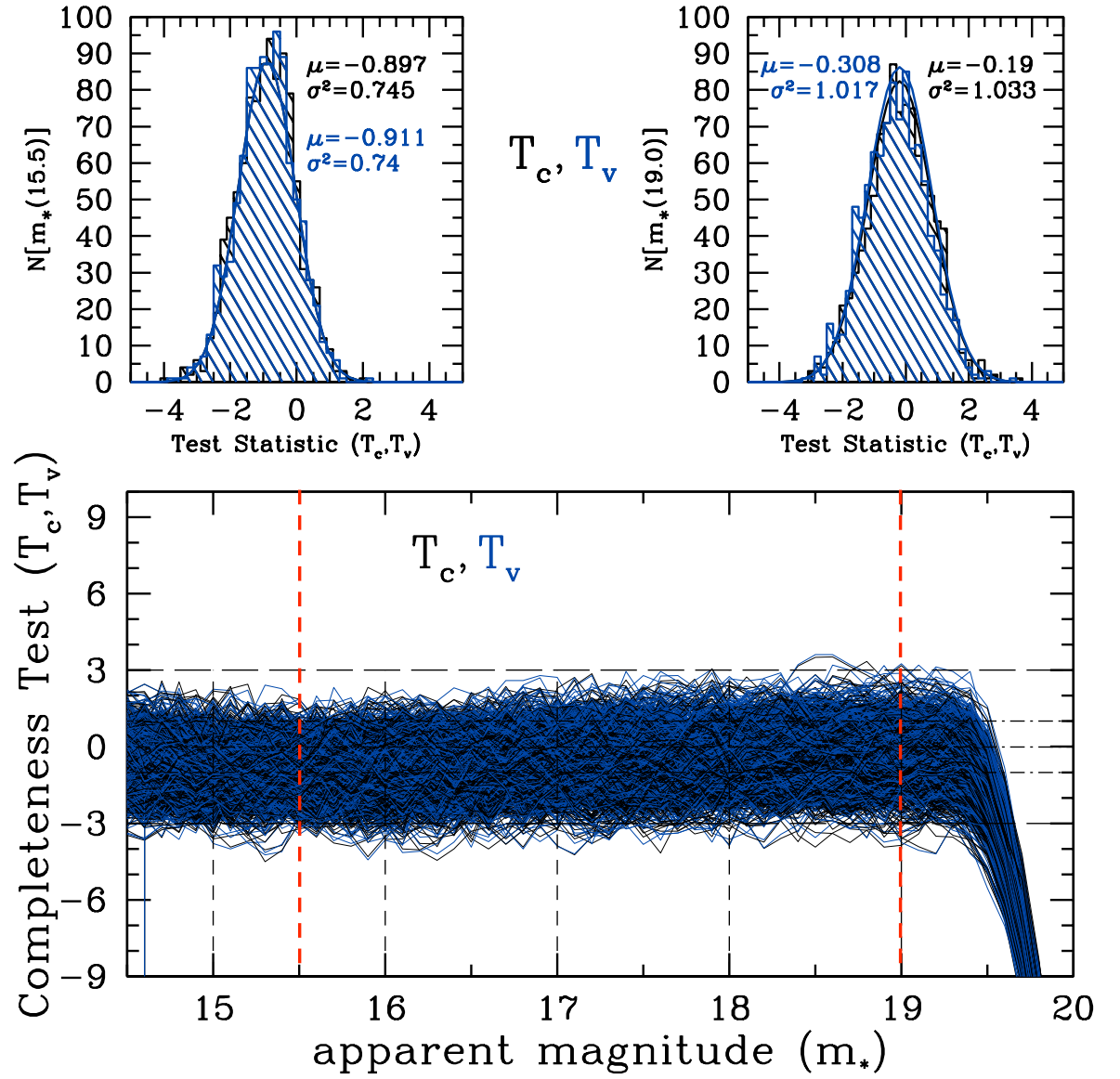


Figure 8.9: T_c and T_v distributions determined from the MGC mocks, applying the R01 method. In this example we have taken slices at $m_* = 15.5$ mag and $m_* = 19.0$ mag indicated by the red dashed line on the bottom panel. Their respective distributions are shown on the top-left and top-right of the figure with Gaussian fits superimposed. On each panel μ and σ^2 refer to the mean and variance of the distributions respectively.

assumption to hold. However, as we will see, we do not necessarily expect the same to apply in the JTH extension under certain conditions. We therefore use our MGC mocks to briefly examine this.

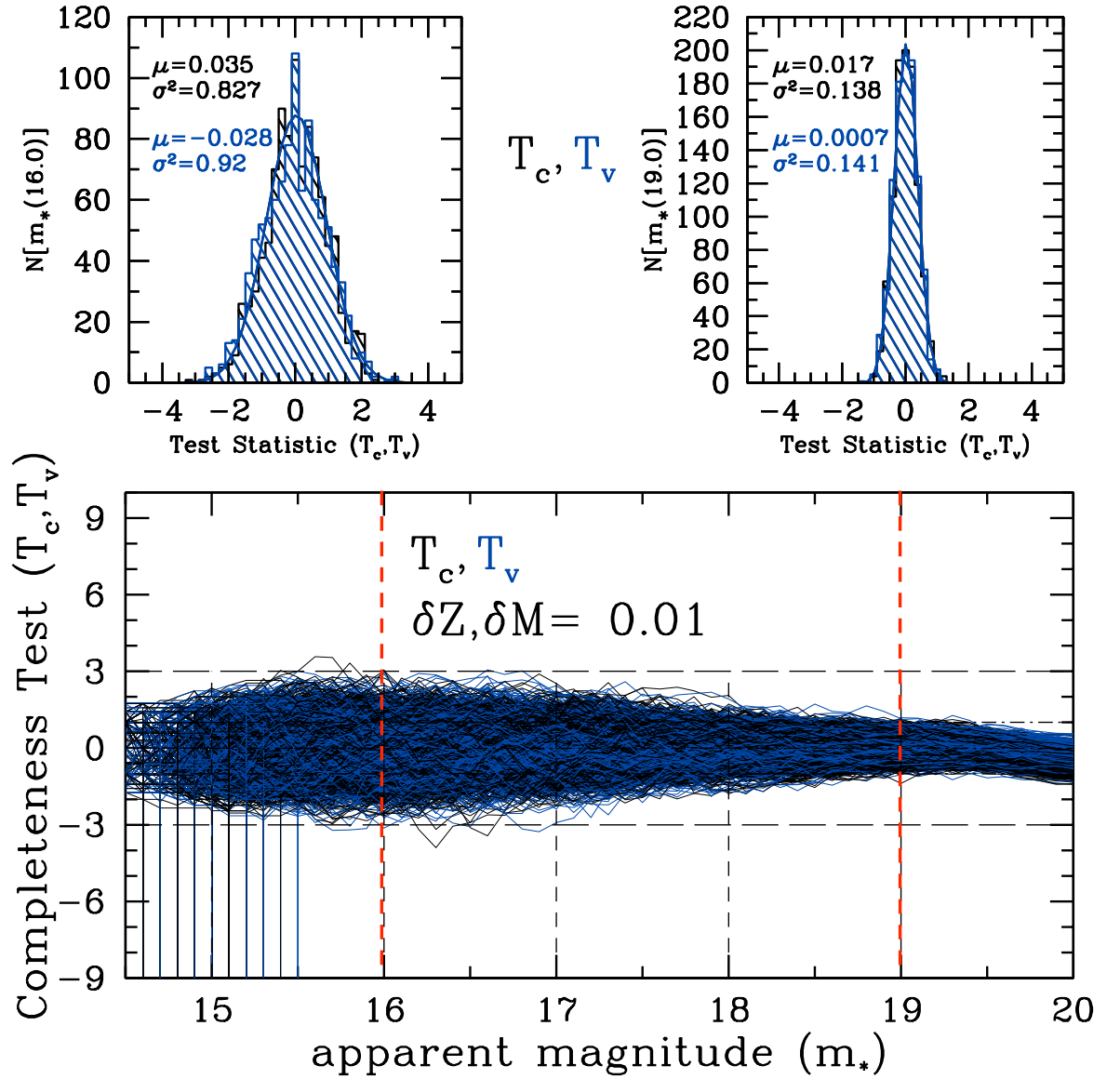


Figure 8.10: T_c and T_v distributions determined from the MGC mocks, applying the JTH method for $\delta Z, \delta M = 0.01$. In this example we have taken slices at $m_* = 16.0$ mag and $m_* = 19.0$ mag indicated by the red dashed line on the bottom panel. Their respective distributions are shown on the top-left and top-right of the figure. On each panel μ and σ^2 refer to the mean and variance of the distributions respectively.

R01: In the bottom panel of Figure 8.9 we have T_c and T_v results for 1000 of our MGC mocks in the G-Frame when applying the R01 method. The red dashed lines represent the slices in m_* that we consider, the first at $m_* = 15.5$ mag and the second at $m_* = 19.0$. The resulting distributions are respectively shown at the top left and right panels in Figure 8.9 along with Gaussian fits superimposed. On each panel μ and

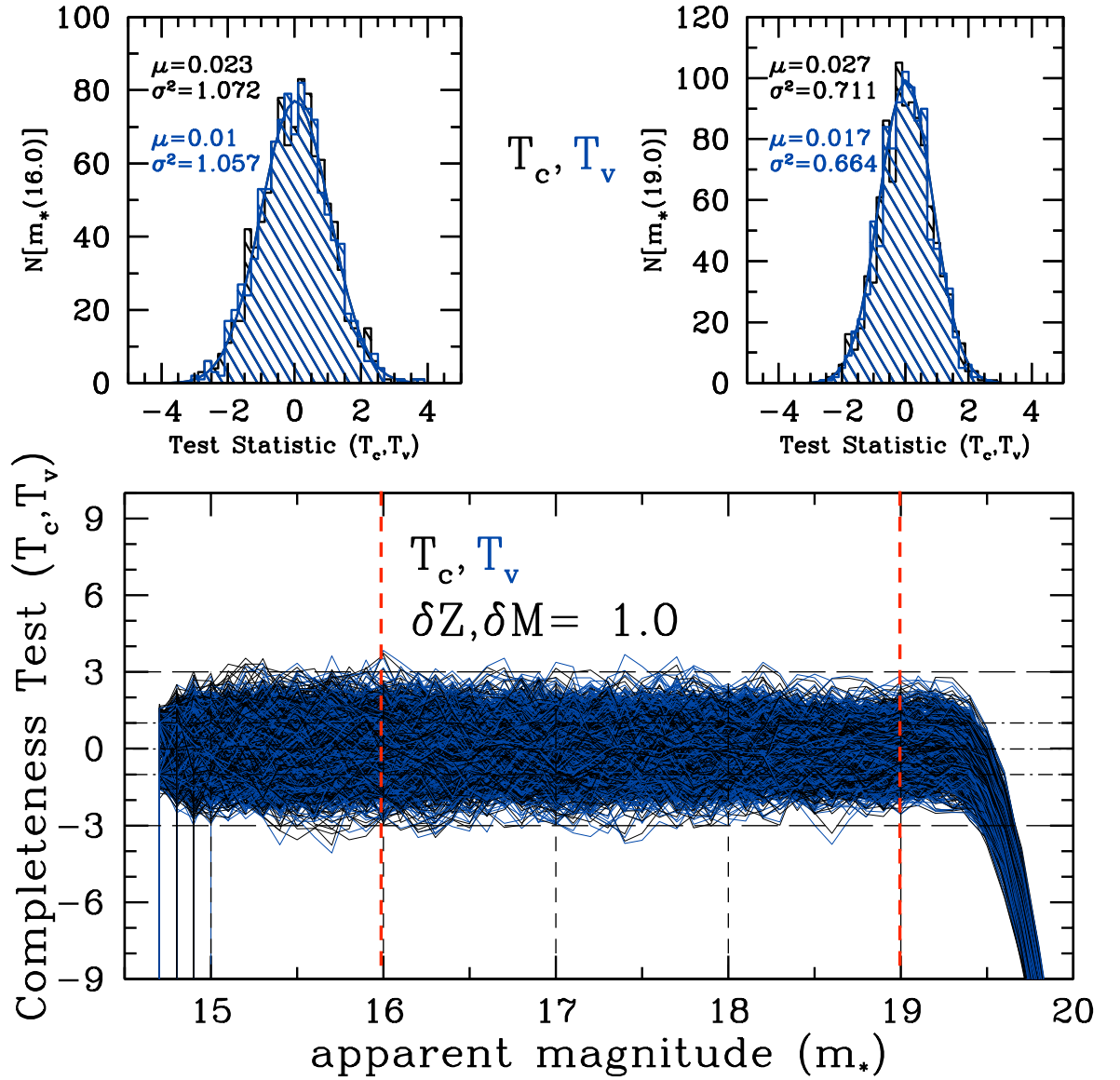


Figure 8.11: T_c and T_v distributions determined from the MGC mocks, applying the JTH method for $\delta Z, \delta M = 1.0$. In this example we have taken slices at $m_* = 15.5$ mag and $m_* = 19.0$ mag indicated by the red dashed line on the bottom panel. Their respective distributions are shown on the top-left and top-right of the figure. On each panel μ and σ^2 refer to the mean and variance of the distributions respectively.

σ^2 refer to the mean and variance of the distributions respectively. It is clear to see at $m_* = 15.5$ mag both T_c and T_v statistics do not have a distribution consistent with mean zero, but in fact have distributions centred around -0.898 and -0.911 respectively. However, they appear to be consistent with being Gaussian with an overall variance of ~ 0.74 . If we check where the complementary red dashed line in the bottom panel

crosses the distribution we can indeed see dip in the overall T_c and T_v distributions. Therefore, this is more a reflection of the mock catalogues as opposed to the underlying properties of the statistics themselves.

Towards the faint magnitude limit at $m_* = 19.0$ we see resulting distributions have a mean $\mu_{T_c}[m_*(19.0)] = -0.190$ and $\mu_{T_v}[m_*(19.0)] = -0.308$ which seems a more consistent result. Moreover both are again consistent with being Gaussian and have variances of $\sigma_{T_c} = 1.033$ and $\sigma_{T_v} = 1.017$.

JTH - $\delta Z, \delta M = 0.01$: We now apply the JTH method to the same mocks and adopt a $\delta Z, \delta M = 0.01$. This value is a very small and has been shown in Chapter 5 to be dominated by shot noise. In Figure 8.10 we can see the two slices at $m_* = 16.0$ mag and $m_* = 19.0$ mag. This time we observe for $m_* = 16.0$ mag mean values of $\mu_{T_c}[m_*(16.0)] = 0.0356$ and $\mu_{T_v}[m_*(16.0)] = -0.028$ with respective variances of $\sigma^2 = 0.827$ and $\sigma^2 = 0.92$. At $m_* = 19.0$ mag we observe a very narrow Gaussian distribution with variances of $\sigma_{T_c}^2 = 0.138$ and $\sigma_{T_v}^2 = 0.141$ and respective means of $\mu[m_*(19.0)] = 0.017$ and $\mu[m_*(19.0)] = 0.0007$. This result confirms once more the effects of shot-noise on our statistics when the widths of δZ and δM are very small.

JTH - $\delta Z, \delta M = 1.0$: Lastly, we increase $\delta Z, \delta M = 1.0$ which we found in previous chapters to be a suitable level to accurately determine the completeness of the data. For fainter magnitudes we now observe results resembling closer to that of the R01 method. This time the mean values are $\mu_{T_c}[m_*(16.0)] = 0.023$, $\mu_{T_v}[m_*(16.0)] = 0.01$ with respective variances $\sigma^2 = 1.072$ and $\sigma^2 = 1.057$. For the faint magnitude slice we find $\mu_{T_c}[m_*(19.0)] = 0.027$, $\mu_{T_v}[m_*(19.0)] = 0.017$ with respective variances $\sigma^2 = 0.711$ and $\sigma^2 = 0.664$.

With this simple analysis, we can see that, overall the mocks validate the underlying properties of T_c and T_v , but have shown us the conditions where they deviate from this - when δZ and δM have values so small that shot noise begins to dominate.

8.5.2 Completeness of 2dFGRS mocks

When creating the 2dFGRS mocks we allowed for a brighter faint apparent magnitude limit in our selection criteria. This resulted in the magnitude distribution being well described by a faint apparent magnitude limit *only*. This is clearly visible in the $M-Z$ distribution shown for mock #78 in the middle-left panel in Figure 8.6. Consequently,

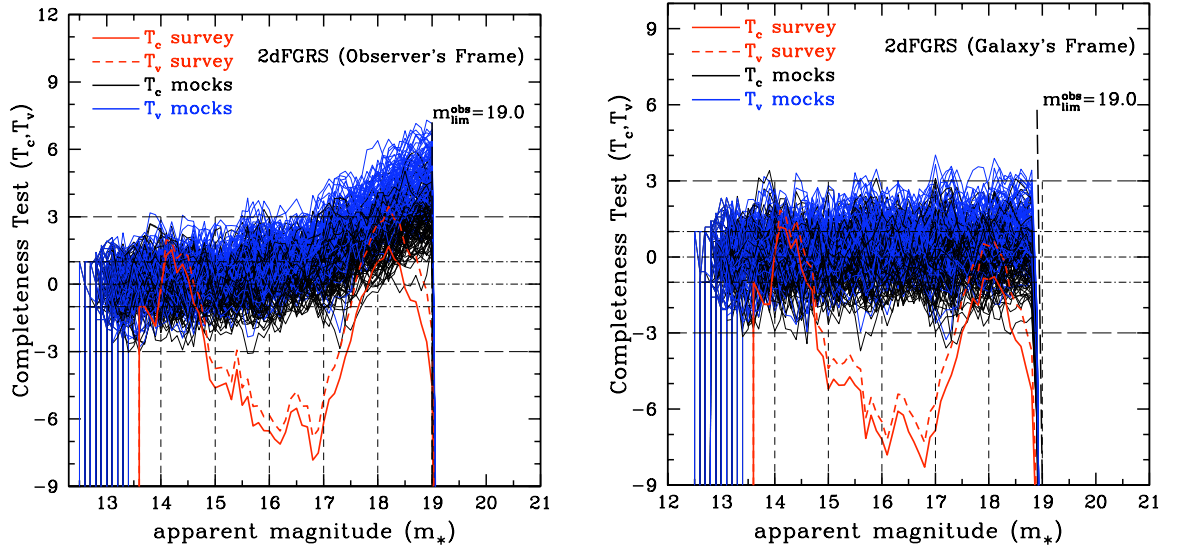


Figure 8.12: T_c and T_v results for 2dFGRS mocks in both the G-Frame and the O-Frame. Here, we apply the R01 methodology to the data which considers a faint apparent magnitude limit only. In the O-Frame on the left-hand panel the same rising trend in T_c and T_v is observed for the mocks as with MGC. Of course, the actual 2dFGRS survey data shown in red, shows a similar trend to that of our original completeness results in § 4.3 on page 4.3.

this has allowed us to apply both the R01 and JTH completeness methods to the mocked data.

If we firstly examine the completeness results for the R01 method in Figure 8.12 then we record a similar result to that of the MGC mocks. In the un-corrected O-Frame the mocks demonstrate a systematic rise in T_c and T_v that begins at around $m_* \sim 16.0$ and continues to rise, just like MGC, to the magnitude limit of the sample, $m_{\text{lim}} = 19.0$ with a peak in $T_v \sim 7.2\sigma$. The survey data indicated by the red lines show the characteristic trend of our original completeness results in § 4.3 on page 4.3.

In the G-Frame shown in the right-hand panel in Figure 8.12, the mock surveys show an overall completeness, with T_c and T_v indicating the true limit at $m_{* \text{lim}} \sim 18.9$ mag due to the curving of the data at the limit.

We now introduce a bright apparent magnitude limit ($m_{\text{lim}}^{\text{b}} = 12.0$) for the mock and real survey data and run the JTH method with a δZ and $\delta M = 0.08$ applied to T_c and T_v respectively (see Figure 8.13). What is immediately obvious in the O-Frame results is the suppression of rising behaviour that was observed in the application of R01 in Figure 8.12. This indicates that for small widths of δZ and δM , any potential effects which would break the separability between M and Z may be masked. However, it should also be noted that in both the O-Frame and the G-Frame of Figure 8.13 the

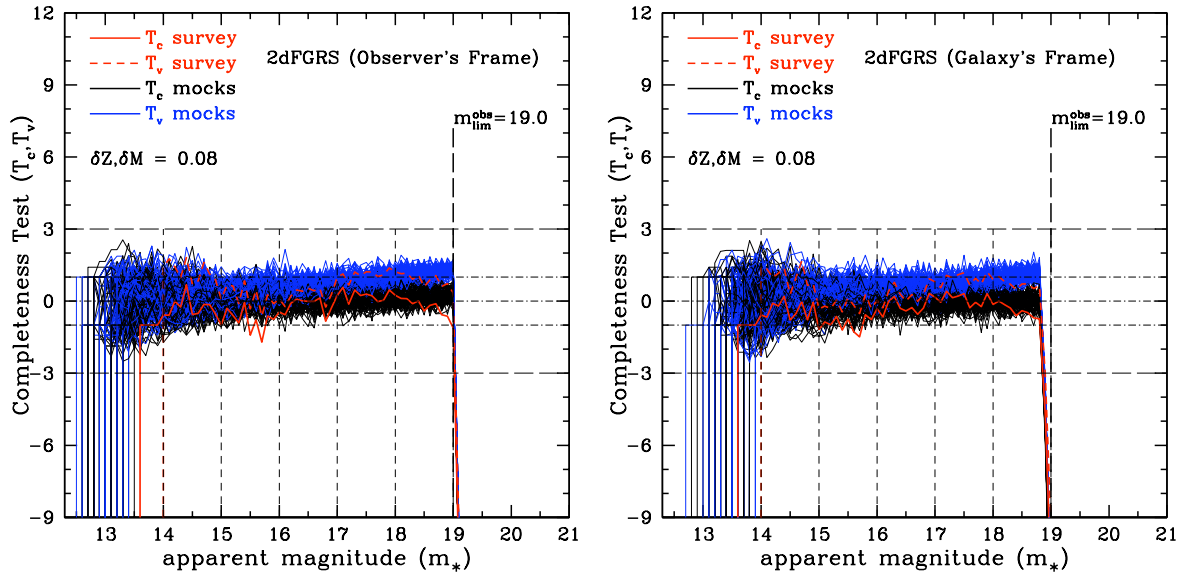


Figure 8.13: T_c and T_v results for 2dFGRS mocks in both the G-Frame and the O-Frame. In this case we apply the JTH methodology to the data which considers both faint and bright apparent magnitude limits.

correct magnitude limit is correctly identified.

8.5.3 Completeness of the SDSS mocks

As we move to our final set of mocks, we note that we should be limited to using the JTH method on SDSS since, as shown on the bottom panels of Figure 8.6, the survey data is defined by a relatively narrow range in apparent magnitude i.e. $14.50 < m < 17.45$, i.e. with sharp cut-off's at both faint and bright apparent magnitude limits. Consequently, applying appropriate widths in δZ and δM when we examine the mocks in the O-Frame, may mask any breaks in separability due to k and/or evolutionary effects. Indeed, in the left-hand plot in Figure 8.14 for the uncorrected data the T_c and T_v statistics for the mocks indicate the M - Z distributions are complete up to the faint magnitude limit indicating that k -corrections have minimal impact on this sample.

In the G-Frame we observe the resulting true limited indicated by T_c and T_v as a due to the curved apparent magnitude limit of the mock sample. In this case, the limit is identified as being $m_{\text{lim}} \sim 17.2$ mag. It also worth noting that the T_c and T_v results for the actual SDSS survey compare very well to the mock T_c and T_v distributions.

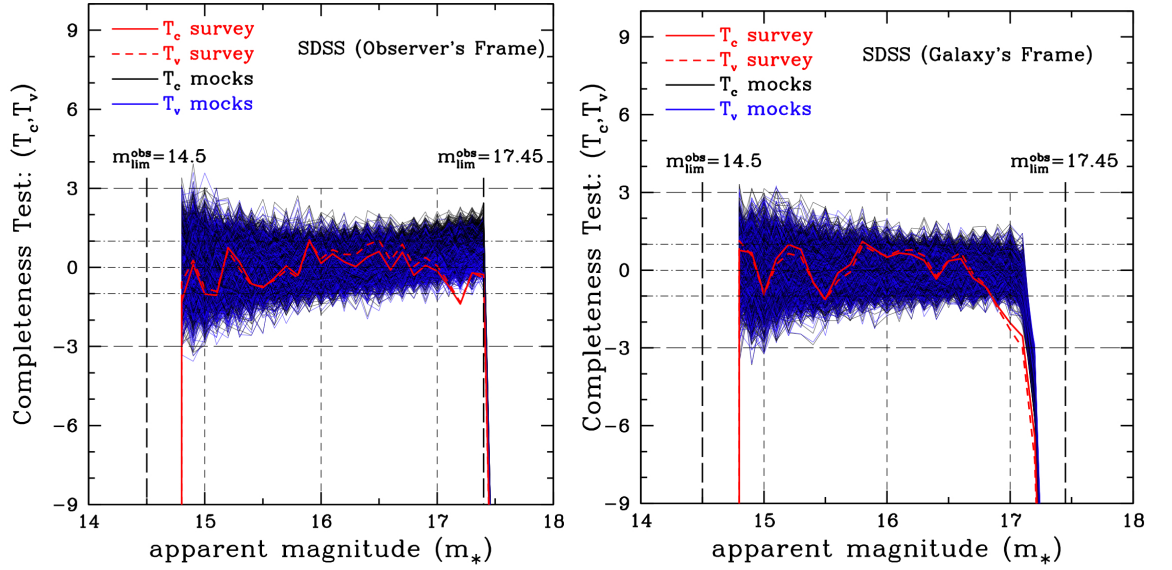


Figure 8.14: T_c and T_v results for 2dFGRS mocks in both the G-Frame and the O-Frame applying the JTH methodology. Since we are using the same sample size as originally defined in Chapter 4, we adopt the same δZ and δM values of 0.2. As was the case with the 2dFGRS mocks, the introduction of δZ and δM appears to mask the break in separability that we would expect to observe in the uncorrected O-Frame. However, in both frames, the correct faint apparent magnitude limits are identified with an average $m_{\text{lim}} \sim 17.2$ mag observed in the corrected G-Frame.

8.6 Conclusions

In this chapter we have examined some of the different approaches one can take to simulate a real galaxy redshift catalogue. We concluded that for our purposes, using an existing observed redshift distribution from a survey was the most efficient and convenient way to proceed as it replicates the observed redshift distribution *exactly* giving us built in clustering. In our examples we used the redshift distributions for MGC, 2dFGRS and SDSS redshift survey samples. We established that the ‘observer’s approach’ allowed us to simulate these galaxy catalogues in two different frames.

The first frame we referred to as, the ‘observer’s frame’, or the O-Frame (not to be confused with the ‘observer’s approach’). In this frame we simulated the real observed apparent and absolute magnitudes of a given survey which represent magnitudes before k - and/or e -corrections have been applied. The second frame was the ‘galaxy’s frame’, or the, G-Frame. This represented apparent and absolute magnitudes that *have* been corrected for k - and/or e -corrections and therefore should be the magnitude of a galaxy at the time of emission, or equivalently at $z = 0$.

We have provided a transparent step by step account of our galaxy selection procedure and compared our mock catalogues to that of the real survey data. We concluded that all our Monte Carlos simulations compared very well real survey in both the G-Frame and O-Frame. Since our 2dFGRS and SDSS survey samples did not represent the exact samples used to derive the survey luminosity function parameters, we required to tweak them slightly so our distributions closely match the observed magnitude distribution.

Once satisfied that we had simulations that represented the three surveys, we performed a completeness analysis using T_c and T_v . We demonstrated for all three survey mocks in the G-Frame, the T_c and T_v statistics systematically dropped at a point that indicated the true magnitude limit brighter than the actual survey limit. In the case of MGC, the survey limit was $m_{\text{lim}} = 20.0$ mag, but T_c and T_v indicated a limit in the range $19.3 \lesssim m_*^{\text{lim}} \lesssim 19.7$. The 2dFGRS mocks had a survey limit of $m_{\text{lim}} = 19.0$ mag but T_c and T_v indicated the magnitude limit to be at $m_*^{\text{lim}} \sim 18.9$ mag. Finally, SDSS whose survey limit was $m_{\text{lim}} = 17.45$ mag showed a resulting T_c and T_v magnitude limit $m_{\text{lim}} \sim 17.2$ mag. The reason for these results lay in the way we selected the galaxies in the G-Frame. By making our final selection of galaxies according to Equations 8.12 and 8.13 which are dependent on the k - and e -correction models, we inevitably introduce a curved apparent magnitude limit. Since our modelling of m_* is such that it is assumed to remain straight, the T_c and T_v statistical results then reflect at which point the curved faint limit begins to dominate.

Conversely, the completeness results in the O-Frame indicated the true apparent magnitude limit to be exactly at the survey limit for all three surveys. Since we convert the selected magnitudes from the G-Frame to the O-Frame by Equation 8.14, this renders the M - Z distribution with straight magnitude limits which resembles the distribution of the raw magnitudes from the real surveys.

Finally, we took a brief look at more fundamental aspect of our statistics that relates to the sampling distribution of T_c and T_v . By their construction they should have a Gaussian sampling distribution with mean zero and variance equal to unity. Using MGC as an example, we compiled T_c and T_v results for a 1000 mock realisations (in the G-Frame) by firstly applying the R01 method followed by the JTH method for a δZ and δM of 0.01 and 1.0. We then sliced through the resulting T_c and T_v distributions at two different m_* values. We showed that for the R01 method and the JTH method for a δZ and δM of 1.0, the underlying properties generally hold. When we applied the JTH method we found that for a small value of δZ and $\delta M=0.01$ we confirmed that for brighter values of m_* the overall sampling distribution for both T_c

and T_v seems consistent with a Gaussian distribution over the $|3\sigma|$ range. As we moved to a faint $m_*=19.0$ mag the distribution was considerably narrower further illustrating the dominance of shot-noise for very small values of δZ and δM . As we increased δZ and δM to 1.0 we observed the distributions broaden as we would expect since we are allowing more galaxies into the each calculation of ζ and τ .

Chapter 9

A New Probe for Evolution

“What I find most disheartening is the thought that somewhere out there our galaxy has been deleted from somebody else’s sample.”

Dr Alexander Boksenberg, 1996

As we move into the final part of the thesis we focus on another major area in characterising the galaxy luminosity function (LF) - evolution. In this chapter we look solely at the role of pure luminosity evolution (PLE) models as applied to galaxy luminosity functions. In Chapter 1 we explored the different effects of evolution on the shape of the LF. Generally, this form of evolution is described parametrically by,

$$e(z) = -2.5\beta \log_{10}(1+z), \quad (9.1)$$

where β is the evolutionary parameter and is usually considered to be galaxy dependent.

In this chapter we will show we can extend our application of robust statistical methods to probe PLE models. We revisit the properties of our random variable, ζ , and introduce a new random variable, χ , based on the observed redshift distribution. We will demonstrate how the ζ - χ distribution derived from a set of magnitudes (drawn from an evolving LF) and redshifts can be used to constrain effectively the value of the evolutionary parameter, β_{true} .

As a controlled testing ground we generate mock galaxy catalogues based on the MGC, 2dFGRS and SDSS data-sets which we have already analysed in Chapter 8 and apply two different techniques intended to estimate the correct value of β . The first approach tests the correlation between our two random variables (ζ, χ), whilst the second approach measures the relative change of entropy between the (ζ, χ). In the

following sections we detail these techniques, and step through our methodology of their application before presenting the results and conclusions.

9.1 Methodology to Probe Evolution

9.1.1 ζ and the new random variable, χ

The cornerstone of our analysis lies in the use of the random variable, ζ . We recall that the definition of ζ for the R01 and JTH case is based on the cumulative luminosity function, denoted as $F(M)$, and respectively given by,

$$\zeta(\text{R01}) = \frac{F(M)}{F[M_{\text{lim}}(Z)]}, \text{ or } \zeta(\text{JTH}) = \frac{F(M) - F[M_{\text{lim}}^b(Z - \delta Z)]}{F[M_{\text{lim}}^f(Z)] - F[M_{\text{lim}}^b(Z - \delta Z)]}. \quad (9.2)$$

For a *complete* sample ζ has, in both cases, the property of being uniformly distributed between $[0,1]$ and independent of the distribution in Z .

We introduce a second variable, χ , that could be used with ζ as a test of independence. This variable is based on the cumulative observed redshift distribution and given by,

$$\chi = \frac{H(z)}{H(z_u)}, \quad (9.3)$$

where,

$$H(z) = \int_z^{z_u} h(z') dz'. \quad (9.4)$$

where z_u is the upper redshift limit. This simple transform means that by definition, χ shares the same property as ζ : its sampling distribution should be uniform on the interval $[0,1]$. For a complete sample, therefore, the sampled points will be uniformly distributed on the ζ - χ plane within a unit square, as illustrated in Figure 9.1. On the other hand; it follows that for a data-set which is drawn from an evolving LF the evolution will have a joint ζ - χ distribution modified by evolution. This will introduce correlations between the proposed variables (see right-hand panel of Figure 9.1). To detect such correlations, we will make use of the two estimators detailed below.

9.1.2 The coefficient of correlation approach

The coefficient of correlation is a well established statistical tool that measures the correlation between two random variables X and Y (or in our case ζ and χ). The

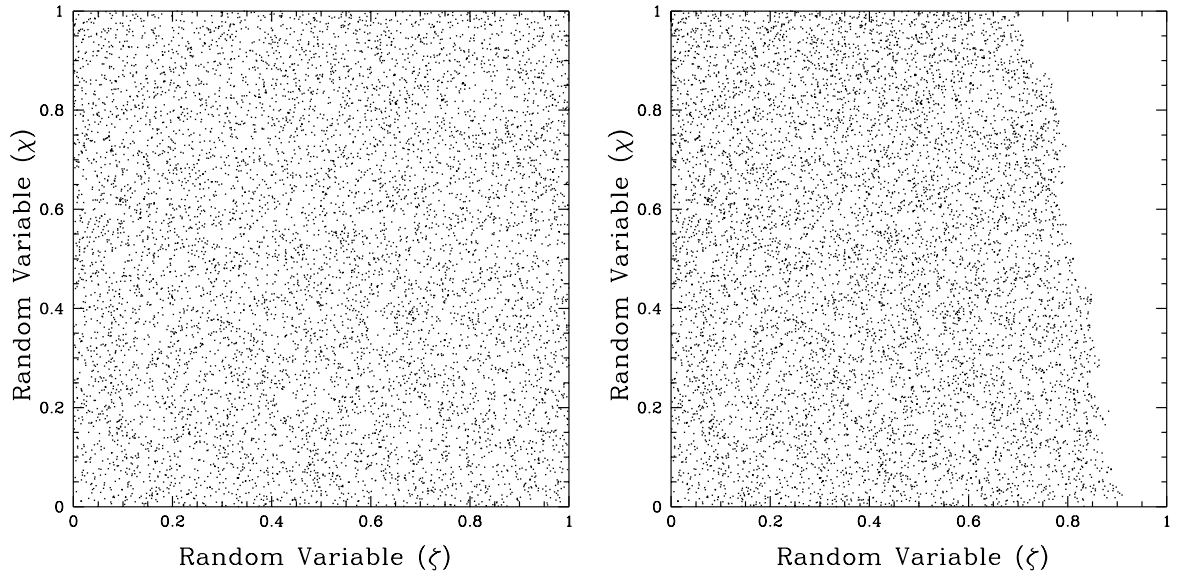


Figure 9.1: Example illustrating a typical (ζ, χ) distribution for a *complete* data-set for an MGC mock catalogue. The left-hand distribution shows ζ and χ estimated at apparent magnitude limit of the survey $m_{\text{lim}} = 20.0$ and appears to be a random uniform distribution. Correlations between ζ and χ are shown on the right-hand panel where ζ and χ have been estimated at $m_* = 20.5$ (beyond the limit of the survey).

coefficient, $\rho(\zeta, \chi)$ is defined as,

$$\rho(\zeta, \chi) = \frac{\text{cov}(\zeta, \chi)}{\sqrt{\text{cov}(\zeta, \zeta)} \sqrt{\text{cov}(\chi, \chi)}}. \quad (9.5)$$

More specifically, ρ is usually estimated from sampled bivariate data using the Pearson product-moment correlation coefficient, and the covariance defined by,

$$\text{cov}(X, Y) = \frac{1}{N_{\text{gal}} - 1} \sum_{i=1}^{N_{\text{gal}}} (X_i - \bar{X})(Y_i - \bar{Y}), \quad (9.6)$$

where N_{gal} is the number of galaxies in the sample and the mean of, \bar{X} and \bar{Y} are simply estimated by,

$$\bar{X} = \frac{1}{N_{\text{gal}}} \sum_{i=1}^{N_{\text{gal}}} X_i, \quad \text{and} \quad \bar{Y} = \frac{1}{N_{\text{gal}}} \sum_{i=1}^{N_{\text{gal}}} Y_i, \quad (9.7)$$

respectively. Therefore, $\rho(\zeta, \chi) \approx 0$ implies no correlation between ζ and χ . This condition will be satisfied if ζ and χ are independent. The implementation of this approach is discussed in § 9.1.4.

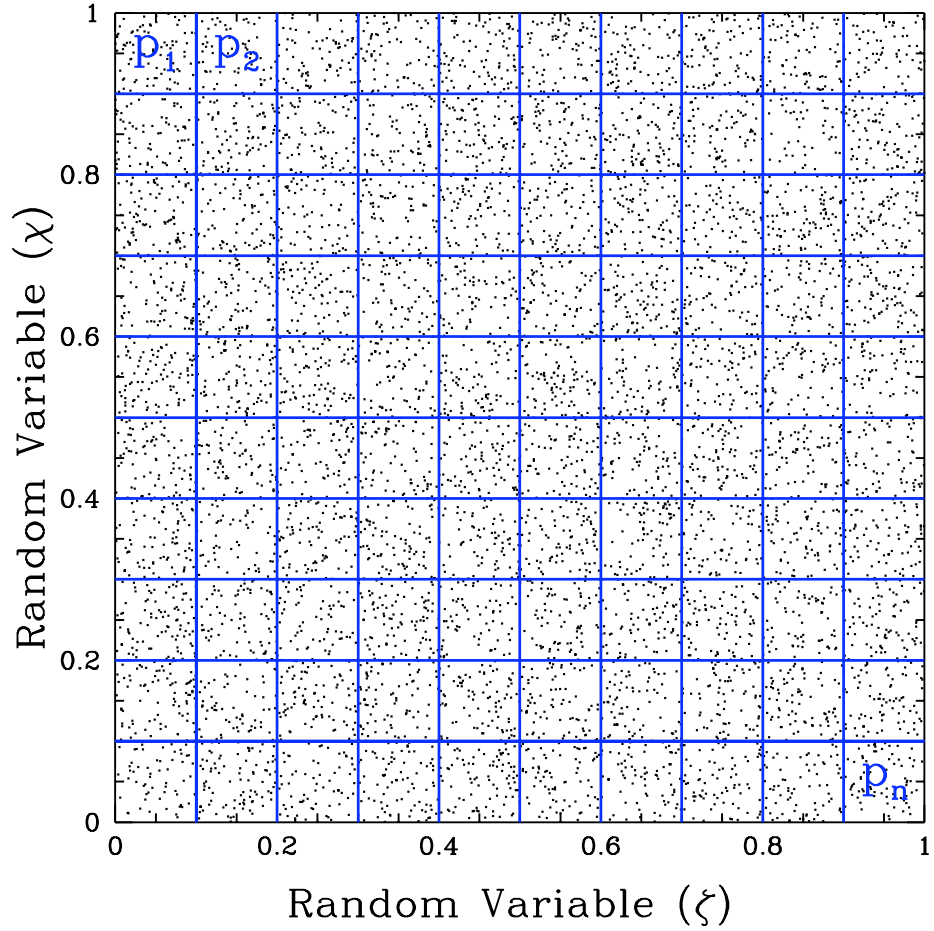


Figure 9.2: Example of measuring the entropy of a typical (ζ, χ) distribution for a *complete* data-set taken at the $m_* = m_{\text{lim}} = 20.0$ of an MGC mock catalogue. We calculate the total entropy of this distribution by imposing a grid with a predetermined mesh size. In this example we have split the grid into 0.1×0.1 mesh. We then count the number of objects contained in each box, p_i , and calculate the relative entropy. We then determine the total entropy by summing each $p_1, p_2, \dots, p_i, \dots, p_n$.

9.1.3 A relative entropy approach

In our second approach we applied a method based on the measured entropy, S , of the ζ - χ distribution. In our case we utilise the relative entropy, also known as the Kullback-Leibler divergence (Kullback and Leibler, 1951), which measures the difference between two probability distributions p and q , where p represents the observed distribution of the data and q represents our theoretical model for that distribution. In Figure 9.2 we illustrate how the entropy is measured for the (ζ, χ) distribution.

We firstly impose a mesh over the distribution with equally spaced cells of a predetermined breadth and height. In the example shown in Figure 9.2 we choose a cell

size of (0.1×0.1) which equates to a total of 100 cells. We will discuss the influence of the mesh size on the final results in § 9.3. For each cell we calculate the probability, p_i from the observed (ζ, χ) distribution given by,

$$p_i = \frac{N(c_{mn})}{N_{gal}}, \quad (9.8)$$

where $N(c_{mn})$ is the number of galaxies within the cell, c , located at (m, n) , and N_{gal} is the total number of galaxies for the whole (ζ, χ) distribution in the sample. We also require to calculate q_i , the theoretical model that represents, in the ζ - χ plane, the ideal uniform distribution of points in the ζ - χ plane such that each cell has the same number of galaxies. This clearly, then, satisfies,

$$q_i = \frac{1}{C_{tot}}, \quad (9.9)$$

where C_{tot} is the total number of cells that make up the mesh. Therefore, for a given (ζ, χ) distribution we can calculate the relative entropy by,

$$S(\hat{\beta}_k) = \sum_{i=1}^n p_i \ln \left(\frac{p_i}{q_i} \right) \quad (9.10)$$

By its definition the value of relative entropy is a convex function of p_i and always non-negative with the property of being equal to zero only if $p_i = q_i$. In reality of course, sample variance within the observed distribution implies that the total relative entropy will always be greater than zero. Therefore the relative entropy of a distribution is ‘maximised’ at its minimum. One could, therefore, redefine S as,

$$S(\hat{\beta}_k) = - \sum_{i=1}^n p_i \ln \left(\frac{p_i}{q_i} \right), \quad (9.11)$$

such that the entropy is maximised at the maximum of $S(\hat{\beta}_k)$.

9.1.4 Implementation

Although the two approaches described in the previous two sub sections can be used to test for correlations in different ways, both are otherwise implemented in exactly the same way. The flowchart in Figure 9.3 summaries the seven steps we carry out to test our evolutionary model, and we now detail those steps below:

1. We create Monte Carlo mock realisations drawn from an evolving luminosity function. This required us to change the procedure slightly that we adopted in

§ 8.3 on page 159. We are applying a pure luminosity evolutionary (PLE) model. This is dependent both on redshift and the evolution parameter, β , which remains fixed thorough the magnitude sampling process. To incorporate this model into the LF we allow the characteristic absolute magnitude, M^* , to evolve, which requires us to re-write Equation 8.10 as,

$$\Phi(M) = \phi_* \exp \{[-c(\alpha + 1)(M - M_{ev}^*)] - \exp[-c(M_{ev}^* - M)]\} \quad (9.12)$$

where, M_{ev}^* , is defined as,

$$\begin{aligned} M_{ev}^*(z) &= M_0^*(z) - E(z, \beta_{true}) \\ &= M_0^*(z) - 2.5\beta \log_{10}(1 + z). \end{aligned} \quad (9.13)$$

Here, M_0^* is the initial value of, M^* , derived for a universal LF from the survey we are simulating. Since M^* is evolving with redshift, we effectively require to create a unique $CDF[M^*(z_i)]$ for each galaxy at a unique redshift z_i .

2. Another way in which these mocks differ from the those derived in Chapter 8 is our final selection of the magnitudes. If our goal in this chapter was to reproduce the universal mocks as in Chapter 8 then the selection conditions as laid out in Equations 8.12 and 8.13 would be recast as,

$$\text{Condition 1} \left\{ \begin{aligned} m_{lim}^f(\text{mock}) &= m_{lim}^f(\text{survey}) - \underbrace{e(z)}_{[1] \text{ LF}(\beta_{true})} - \underbrace{k(z)}_{[2] \text{ from survey}} \\ m_{lim}^b(\text{mock}) &= m_{lim}^b(\text{survey}) - \underbrace{e(z)}_{[1] \text{ LF}(\beta_{true})} - \underbrace{k(z)}_{[2] \text{ from survey}} \end{aligned} \right. \quad (9.14)$$

$$\text{Condition 2} \left\{ \begin{aligned} M_{lim}^f(\text{mock}) &= m_{lim}^f(\text{mock}) - 5 \log_{10}[d_L(z_i)] - 25 \\ M_{lim}^b(\text{mock}) &= m_{lim}^b(\text{mock}) - 5 \log_{10}[d_L(z_i)] - 25 \end{aligned} \right. \quad (9.15)$$

where, in Condition 1, term-[1], $e(z)$ is the evolution model as applied to the luminosity function (LF) in Equation 9.13 and term-[2] represents the k -corrections originally derived from the real survey. Therefore, to simplify our testing we have decided to set $k(z) = 0$ and therefore exclude term-[2] from the final magnitude selection process.

3. We have at this stage a mock that represents the distribution of apparent and absolute magnitudes drawn from an evolved LF. Therefore, the next step is to correct the magnitudes starting with a minimum trial value of $\hat{\beta}_{\min}$.
4. For each corrected sample data-set at $\hat{\beta}_k$, we then estimate our random variables ζ and χ for a trial magnitude limit m_* equivalent to the faintest apparent magnitude, m of the corrected data-set.
5. At this point we feed the resulting (ζ, χ) distributions into our ρ and entropy estimators. For both estimators we apply incremental values $\hat{\beta}$ (typically 0.1) in the range $\hat{\beta}_{\min} < \hat{\beta} < \hat{\beta}_{\max}$ that encapsulates the true β_{true} value used to create the mock. Therefore, as each successive $\hat{\beta}$ is implemented, we expect the value of the ρ estimator to cross zero for a value of $\hat{\beta}$ that corresponds to the β_{true} applied to the LF of the mock survey. Similarly, the relative entropy estimator should minimise for a value of $\hat{\beta}$ that corresponds to the β_{true} . These two scenarios are illustrated in Figure 9.3.
6. Stages 2 to 4 are repeated for the next trial value of $\hat{\beta}$.
7. Stages 1 to 5 are now repeated for a new random seed in the LF magnitude selection to determine results for a new realisation.

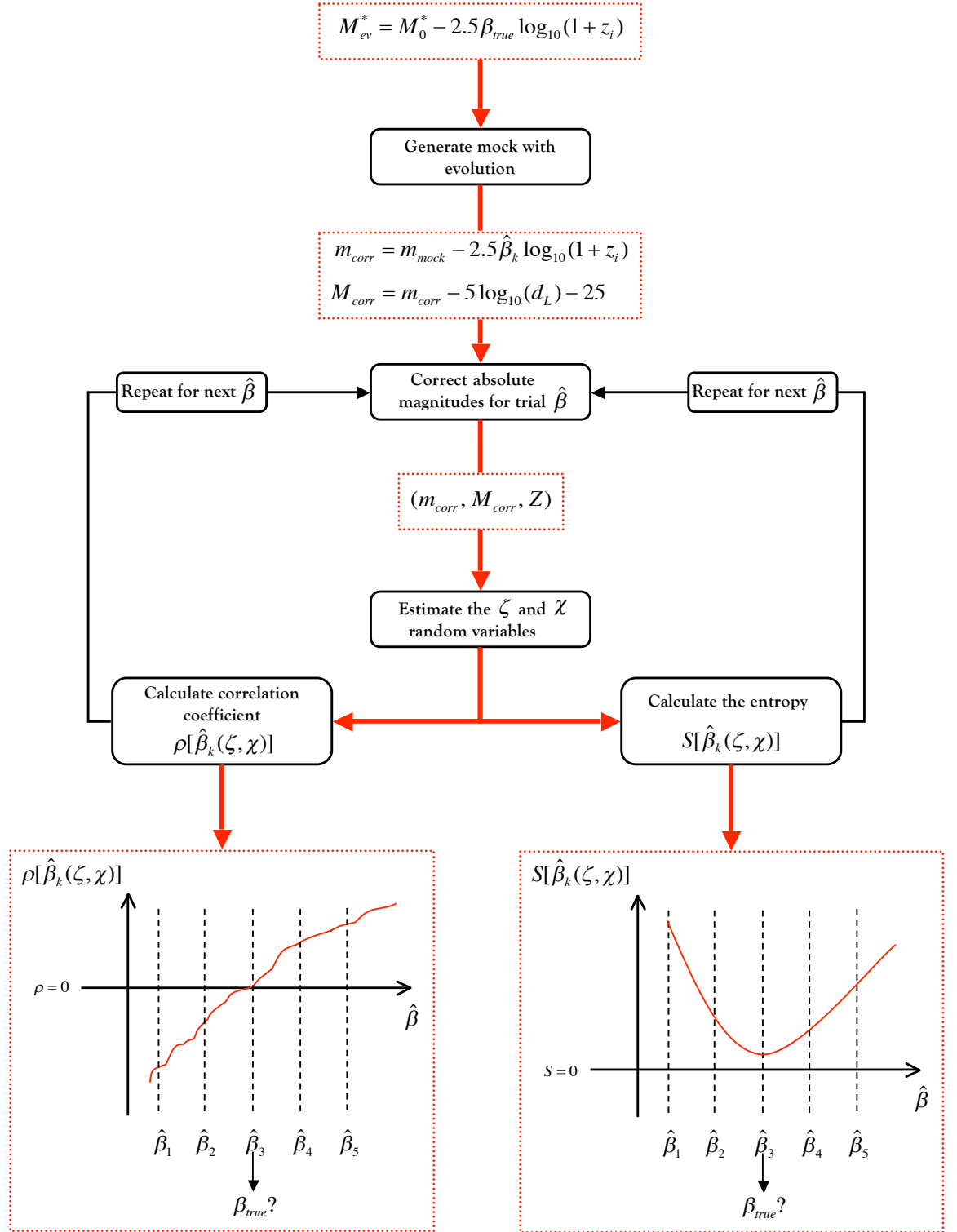


Figure 9.3: Flow diagram summarising the implementation of our test for evolution. The two plots illustrate the kind of results we expect when we apply either the correlation coefficient (left) or the relative entropy approach (right). For the correct value of $\hat{\beta}$ we expect the corresponding ρ value to equal zero (within a given statistical uncertainty). In the case of the relative entropy, we expect the value of S to be minimised at the correct value of $\hat{\beta}$.

9.2 Results Part I - Coefficient of Correlation

For the purposes of applying both the coefficient of correlation method and the relative entropy method, we have decided to use the MGC and SDSS mock catalogues.

9.2.1 MGC - Mocks (R01)

Throughout all the work carried out during this thesis the MGC survey has stood apart from the others that we have analysed, as we have applied both the R01 and JTH method and achieved consistent completeness results which imply the survey is complete and well described with a faint limit only.

Using our MGC mocks to apply the coefficient of correlation approach allows us, therefore, to consider generating (ζ, χ) distributions using either the R01 or the JTH method. For our initial analysis we shall apply R01 to observe the impact where we have no size restriction on the areas S_1 and S_2 (see Figure 3.2 on page 58). We have generated mock catalogues that have no evolution (i.e. $\beta_{true} = 0$) and those drawn from a LF with $\beta_{true} = 1.0, 2.0$ and 3.0 . In Figure 9.4 we present the results for each of these values for a range in trial $-10.0 \leq \hat{\beta} \leq 10.0$ in increments of $\hat{\beta}_{inc} = 0.2$. For an initial analysis we have applied 50 mock realisations in to each β_{true} value in the panels of Figure 9.4.

The top left panel of Figure 9.4 shows the results for $\beta_{true} = 0$. For larger negative values of $\hat{\beta}_k \lesssim -1.0$, the coefficient, $\rho(\zeta, \chi)$ indicates positive-correlation between ζ and χ . Then, in general, we observe that as $\hat{\beta}_k$ approaches zero (i.e the true value, β_{true} used in generating the mocks), $\rho(\zeta, \chi)$ also tends to zero implying no correlation. For the relatively small number of mock realisations that we are considering, there is an obvious spread in the values of beta around $\hat{\beta} \sim 0$ for where $\rho(\zeta, \chi) = 0$ indicates the estimated *true* value of β . The correlation coefficient vanishing (absence of evolution) is obtained for $\hat{\beta}$ in the range $-0.9 \lesssim \hat{\beta}_k \lesssim 0.5$. We can use the spread in $\hat{\beta}(\rho = 0)$ to calculate the uncertainty in $\hat{\beta}$. This is performed more rigorously § 9.2.2.1. For now we are concerned only with the general behaviour of the ρ estimator with $\hat{\beta}$. Finally in this plot we see that as $\hat{\beta}_k$ moves to increasingly larger positive values ρ now indicates strong anti-correlations between ζ and χ and therefore these values of $\hat{\beta}$ can be excluded.

The remaining three panels show the results for positive values of the evolution parameter. If we look at the top-right panel, for $\beta_{true} = 1.0$, we observe a similar spread in $\hat{\beta}$ values for $\rho = 0$ that are this time roughly distributed around $\hat{\beta} = 1.0$ with a range $0.8 \lesssim \hat{\beta}(\rho = 0) \lesssim 1.5$. For $\beta_{true} = 2.0$ and 3.0 the resulting range in $\hat{\beta}$ for $\rho = 0$ are shown to be $1.1 \lesssim \hat{\beta}(\rho = 0) \lesssim 2.4$ and $2.3 \lesssim \hat{\beta}(\rho = 0) \lesssim 3.4$ respectively.

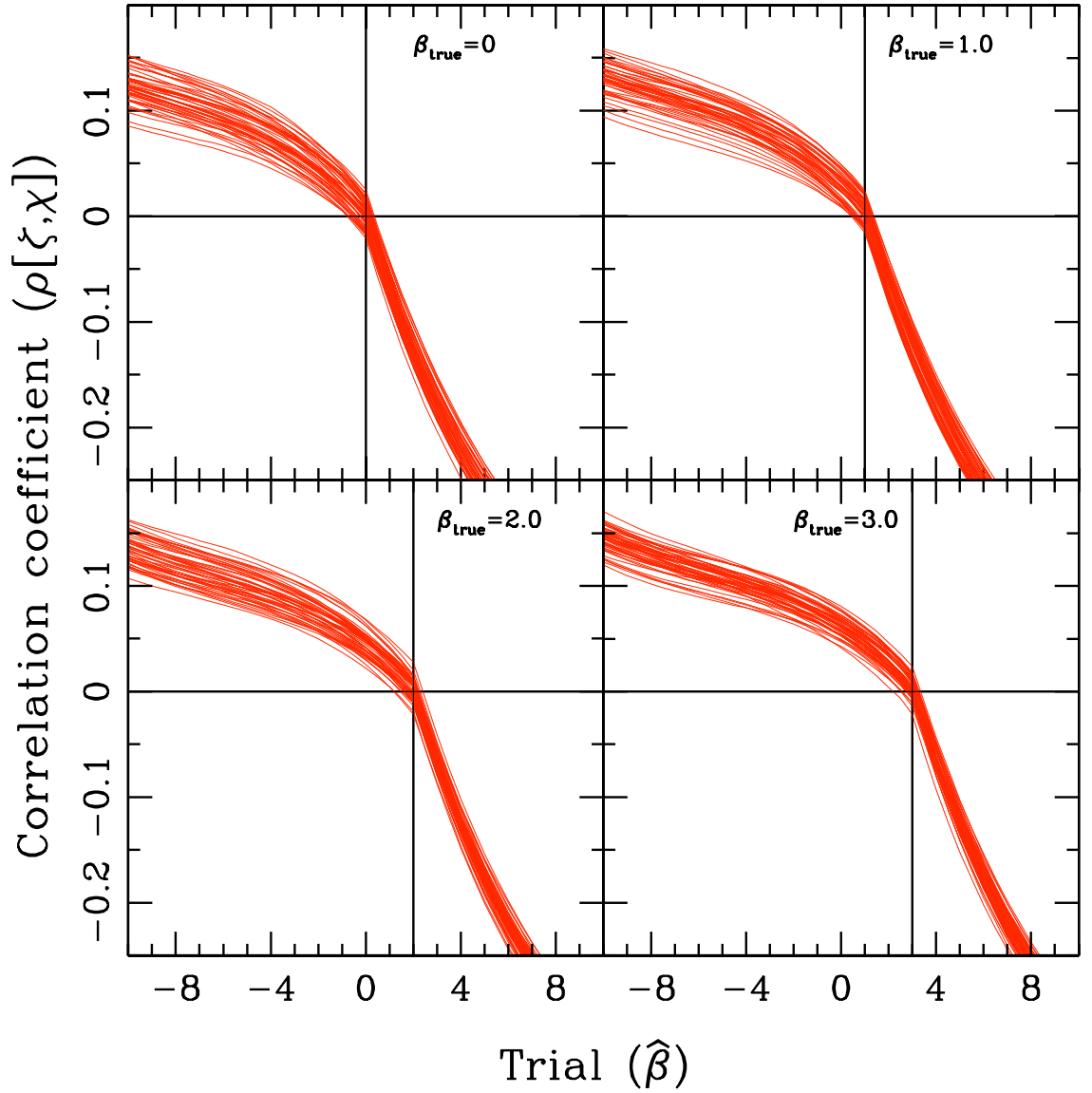


Figure 9.4: Trial $\hat{\beta}$ vs the correlation coefficient ρ for a series of MGC mock catalogues using the R01 method for estimating ζ and χ . For $\beta_{\text{true}} = 0$ we find a resulting range $-0.9 \lesssim \hat{\beta}(\rho = 0) \lesssim 0.5$. The top right panel indicates a range $0.8 \lesssim \hat{\beta}(\rho = 0) \lesssim 1.5$ for $\beta_{\text{true}} = 1.0$ at $\rho = 0$. For $\beta_{\text{true}} = 2.0$ (bottom left) and 3.0 (bottom right) we observe a respective range of $1.1 \lesssim \hat{\beta}(\rho = 0) \lesssim 2.4$ and $2.3 \lesssim \hat{\beta}(\rho = 0) \lesssim 3.4$.

9.2.2 MGC- Mocks (JTH)

We now consider the case where MGC has a bright limit and apply the JTH method to determine the ζ - χ distribution for each trial value of $\hat{\beta}$. In each case we adopt a bright limit corresponding to the brightest galaxy in each $\hat{\beta}_k$ corrected sample. As we demonstrated in Chapter 6, values of δZ that are too small result in shot-noise

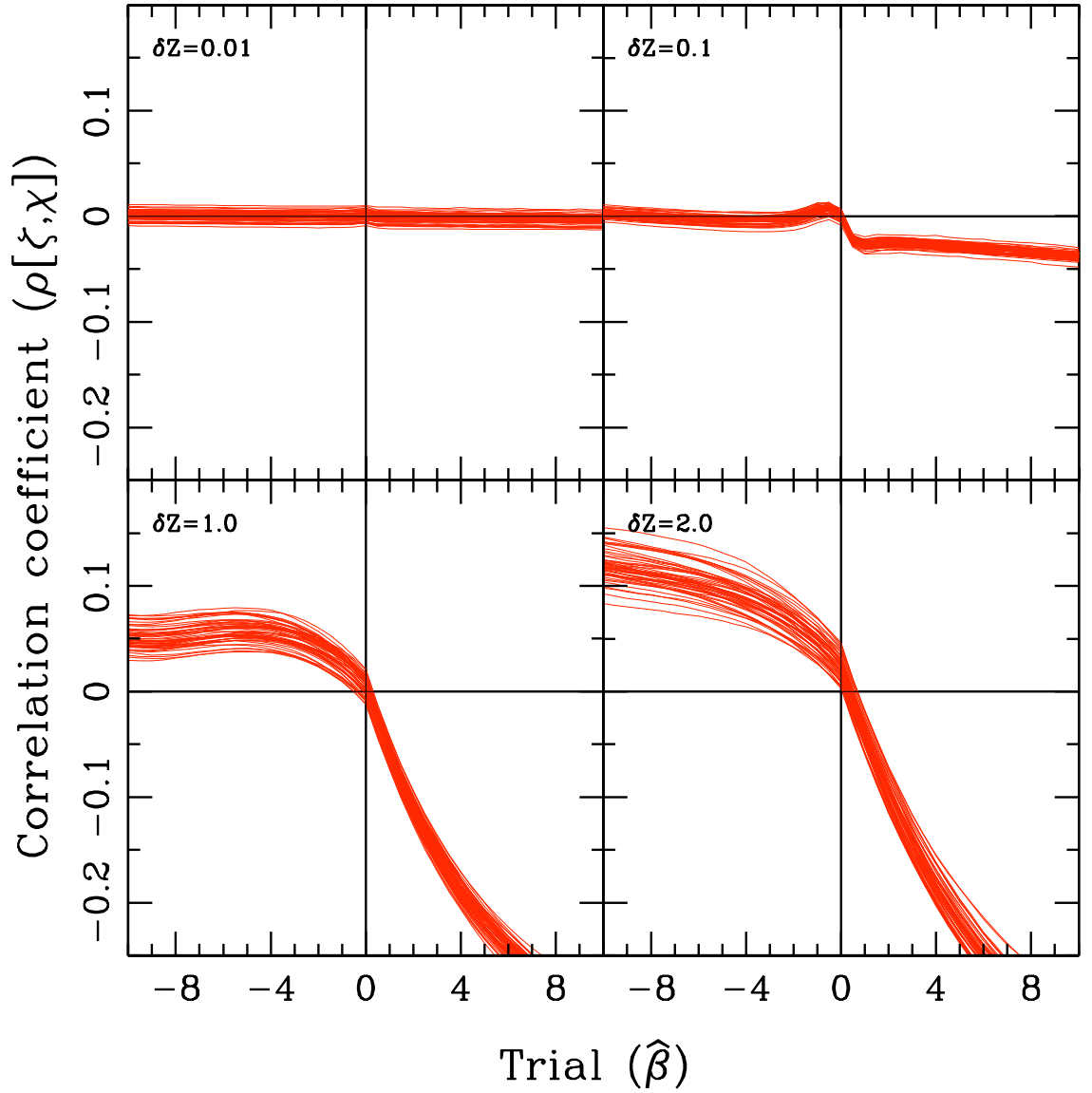


Figure 9.5: Plot showing how varying sizes of δZ effect the ability of the ρ estimator to constrain β_{true} . In each panel we have generated 50 mock realisations. In the top-left panel we have fixed $\delta Z = 0.01$ which results in the $\hat{\beta}$ being distributed about $\rho = 0$ for the entire range of $\hat{\beta}$. This is another indicator of how shot noise can dominate over the signal. Since δZ is so small the resulting number of galaxies in S_1 and S_2 regions is too small to notice any significant change in the ζ estimation for each value of $\hat{\beta}$. Therefore, ρ will approximately equal zero for any $\hat{\beta}$. We observe that a $\delta Z = 1.0$ shown in the bottom left panel, is now large enough for ρ to successfully constrain $\hat{\beta}$.

dominating in the T_c statistics. Therefore, as m_* passes beyond the survey limit, T_c flat-lines since there remains constant small number of galaxies in the resulting S_1 and S_2 regions. Therefore, when applying the JTH method to the ρ estimator we must be

careful in choosing a suitable value of δZ for estimating ζ for the corrected magnitudes according to the trial $\hat{\beta}$.

In Figure 9.5 we have set $\beta_{true} = 0$ and applied the ρ estimator for $\delta Z = 0.01, 0.1, 1.0$ and 2.0 to 50 mock realisations. We show that for $\delta Z = 0.01$ and 0.1 the ρ estimator suffers from the same shot noise issue as with T_c and T_c (see Chapter 6) which is not unsurprising. It is not until that $\delta Z = 1.0$ (bottom-left panel) that we observe ρ successfully constraining β_{true} . Notice also that for $\delta Z = 2.0$, the overall distribution in $\rho(\hat{\beta}) = 0$ appears to systematically shift to a range $-0.8 \lesssim \hat{\beta} \lesssim 0.0$. The most probable reason for this is when δZ becomes sufficiently large, galaxies near the our imposed *bright* apparent magnitude limit can no longer be sampled and as such these galaxies are dropped from the calculation. Therefore, the total number of galaxies in the survey that are used in the ζ calculation is reduced and as a result is now no longer representative of the original mock data-set. Therefore, this implies that ρ in this special case will not be able to recover the original true value of β .

9.2.2.1 Error Analysis

Ultimately, we want to assess the distribution of $\hat{\beta}$ at the point where the correlation coefficient is zero. In this section we have generated more bootstrap (BS) mock realisations in order to accurately assess the distribution and therefore the uncertainty in $\hat{\beta}$ about $\rho = 0$. In order to maximise computing time and give an accurate error analysis we have compared the difference in the $\hat{\beta}(\rho = 0)$ distributions for 200 and 1000 bootstraps as shown in Figure 9.6. Also, for each curve on the $\hat{\beta}$ - ρ plane, we have used interpolation at intervals of $\hat{\beta} = 0.05$ to estimate each value of $\hat{\beta}$ at $\rho = 0$. In the example in Figure 9.6 we have applied the JTH method to the MGC mock data with a $\delta Z = 0.5$. The blue distribution represents 200 BS's, whilst the red represents 1000 BS's. For each case we indicated the mean, $\bar{\beta}$, shown as vertical dashed lines as well as the 68% confidence interval (CI) indicated as solid vertical lines. which further shows a minor change on the right hand error.

For the 200 BS level we determined a $\bar{\beta} = 0.017 \pm 0.001$ and a corresponding CI of $[-0.260, 0.146]$. When we compare this for 1000 BS's we observe very little change with a $\bar{\beta} = 0.017 \pm 0.0001$ and a CI of $[-0.2311, 0.1544]$. Since there is not a significant difference in the distributions we concluded that 200 BS's was enough to make a reasonable assessment of the data. If we now turn to Figure 9.7 we show for each β_{true} value, the 200 BS curves for ρ versus $\hat{\beta}$ applying the R01 method (i.e. assuming a faint apparent magnitude limit) over narrower $\hat{\beta}$ range than in Figure 9.4.

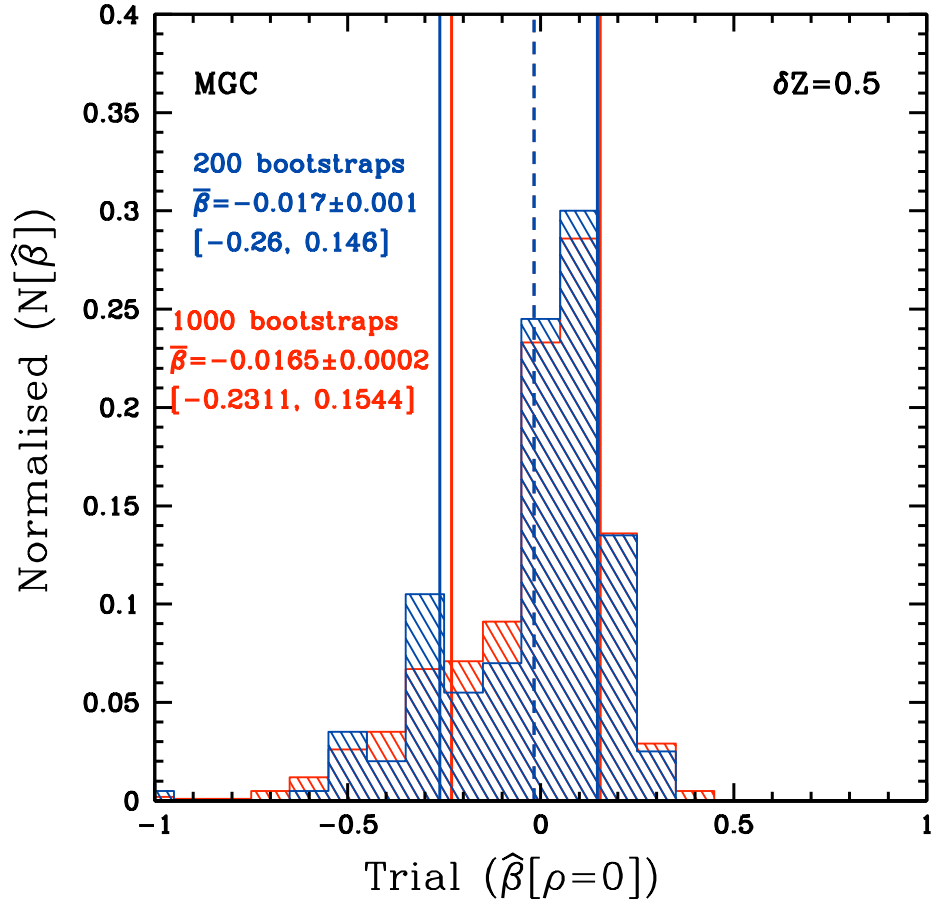


Figure 9.6: MGC $\hat{\beta}$ distribution at $\rho = 0$ for 200 (shown in blue) and 1000 (shown in red) bootstraps. Both distributions were obtained by applying the JTH method to each mock for a $\delta Z = 0.5$. The blue and red vertical lines represent the 68% confidence interval (CI) for the respective 200 and 1000 bootstraps. We can see that increasing the number of bootstraps from 200 to 1000 only marginally affects the $\hat{\beta}$ distribution and consequently the errors at 68% CI.

We observe that for each panel in Figure 9.7 the overall distribution is roughly centred about the expected true β value. However, in each case there is definite skewness to the right of $\hat{\beta}(\rho = 0)$ which is more clearly seen in Figure 9.8. Here, we have plotted the corresponding $\hat{\beta}(\rho = 0)$ distributions for each β_{true} . As with Figure 9.6 we have indicated the 68% CI as blue vertical lines and in blue parenthesis. We have also included the mean of the distribution as a red vertical dashed line with the corresponding value, $\bar{\beta}$, also in red. Although, in each panel we can see that the mean of the $\hat{\beta}$ distribution is consistent with each β_{true} value, each distribution is slightly skewed. The reason for this is best illustrated by considering the (ζ, χ) distributions

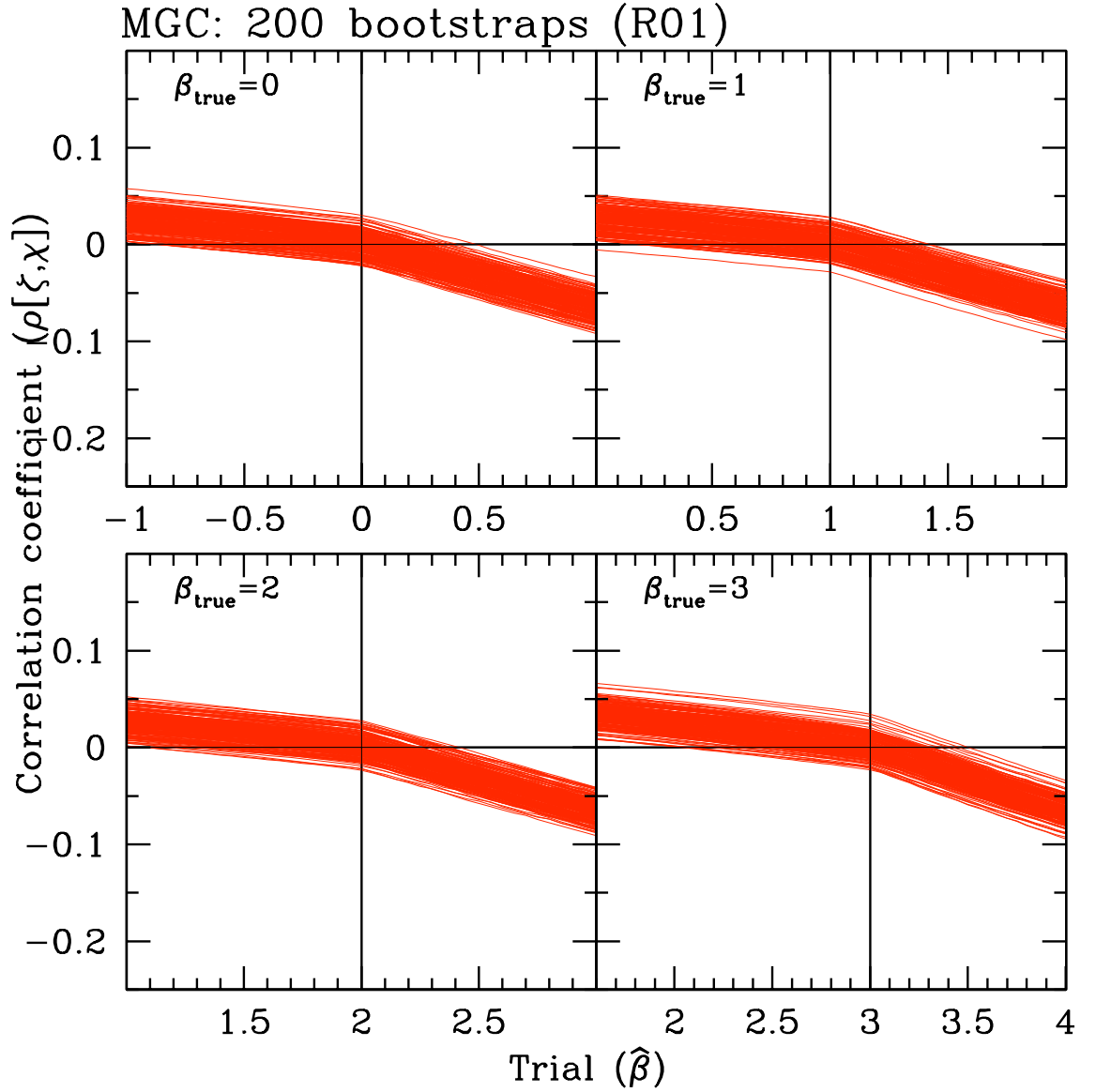


Figure 9.7: 200 MGC bootstraps applying the R01 method for four values of $\beta_{true} = 0.0, 1.0, 2.0$ and 3.0 . In each case we have plotted the resulting ρ versus $\hat{\beta}$ with each red line representing each bootstrap. It is clear to see that for each β_{true} the overall distribution at $\rho = 0.0$ is skewed. Figure 9.8 better shows the $\hat{\beta}$ distribution at $\rho = 0.0$.

resulting from selected $\hat{\beta}$ corrections. In Figure 9.9 we consider one of the MGC mocks with magnitudes that have been sampled with no evolution (i.e. $\beta_{true} = 0$). We have then corrected the magnitudes with trial $\hat{\beta} = \beta_{corr} = -3.0, -10.0$ (left panel) and $\hat{\beta} = \beta_{corr} = +3.0$ and $+10.0$ (right panel) and plotted the resulting (ζ, χ) distributions. In the left panel we can see that from a $\beta_{corr} = -3.0$ (shown in red) to $\beta_{corr} = -10.0$ (shown in black), the distribution becomes extremely anti-correlated. This explains

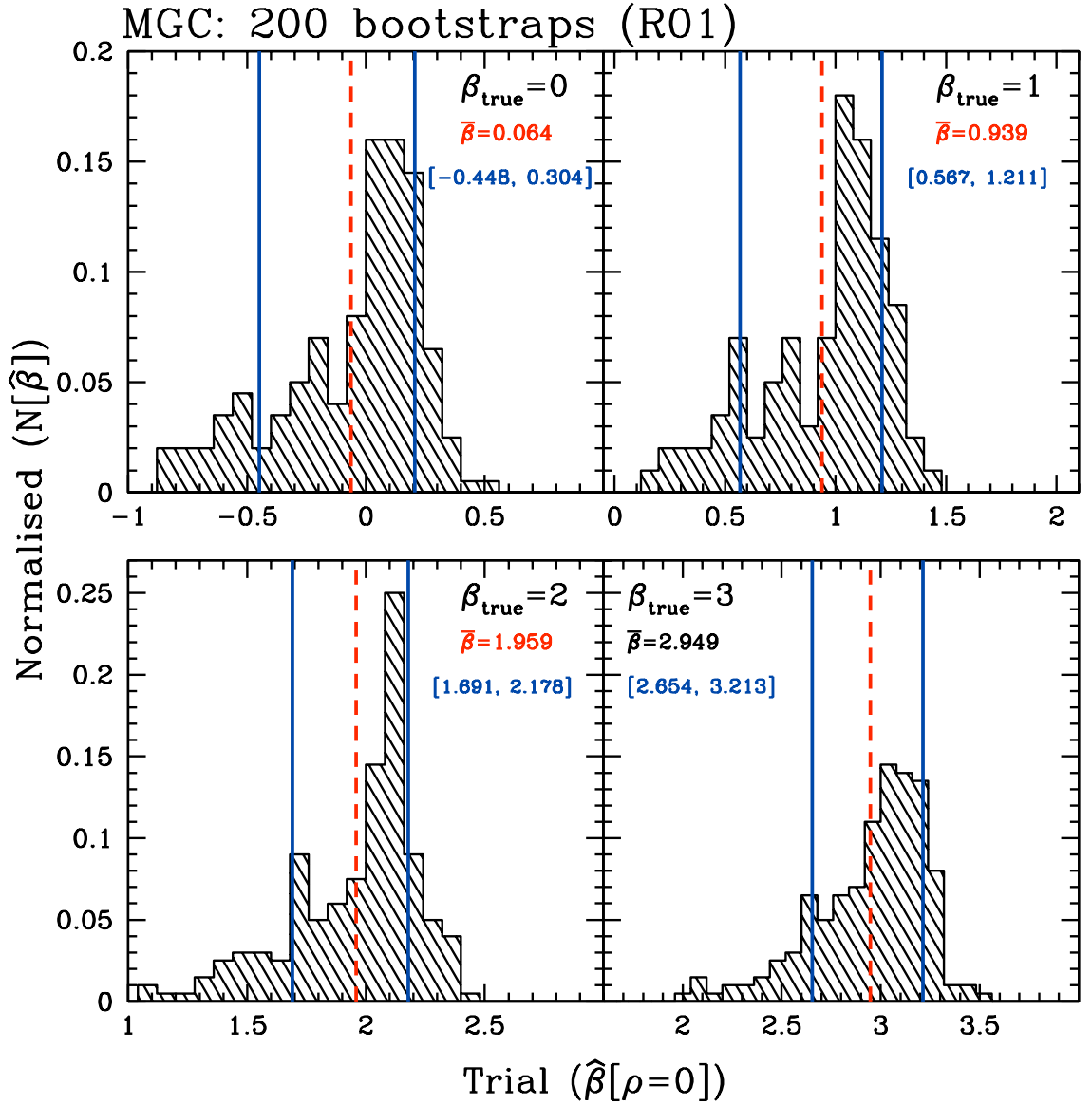


Figure 9.8: MGC bootstrap $\hat{\beta}$ distribution at $\rho = 0.0$. In each panel the red dashed line represents the mean of the distribution, with the corresponding value, $\bar{\beta}$, in red. The blue vertical lines show the 68% confidence interval with the corresponding values in blue parenthesis.

the steepness the $(\hat{\beta}, \rho)$ curves observed in Figure 9.4. However, as we move to positive values of $\hat{\beta}$ (i.e. $\hat{\beta} > \beta_{\text{true}}$) we observe that although the (ζ, χ) distribution is now positively correlated, there is only a marginal change in the distribution $\beta_{\text{corr}} = +3.0$ to $+10.0$. Looking at the corresponding M - Z for the same MGC mock in Figure 9.10 we have plotted the distributions for $\beta_{\text{corr}} = -10.0$ (blue), $\beta_{\text{corr}} = \beta_{\text{true}} = 0$ (black) and $\beta_{\text{corr}} = +10.0$ (green). Also shown are the respective trial limiting apparent magnitudes, $m_* = 19.8, 20.0$ and 21.8 mag, which are determined from the faintest

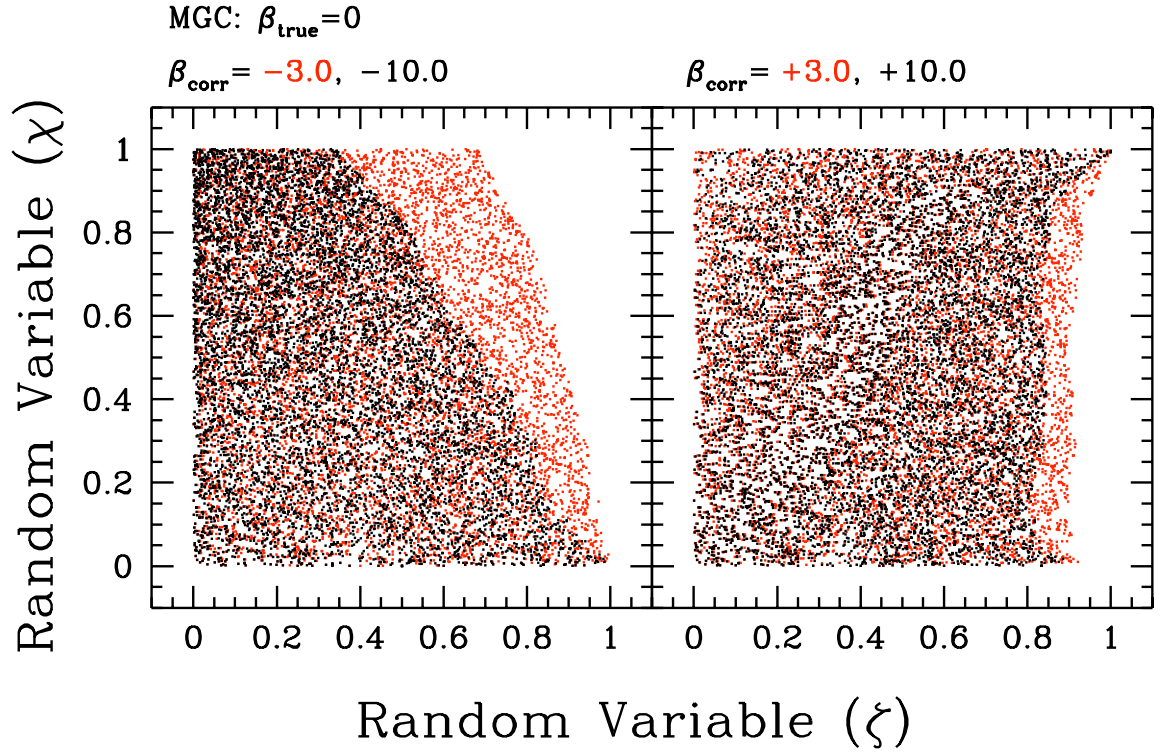


Figure 9.9: ζ - χ distribution for one of the MGC mocks. These plots demonstrate how we obtain anti-correlations and positive correlations for respective $\hat{\beta}$ values which are both smaller and larger than the β_{true} . In this example $\beta_{\text{true}} = 0$ and the left hand panel shows the distribution for $\beta_{\text{corr}} = \hat{\beta} = -3.0$ (shown in red) and $\beta_{\text{corr}} = \hat{\beta} = -10.0$ (shown in black) where we observe anti-correlation. Conversely, the right hand panel shows $\beta_{\text{corr}} = \hat{\beta} = +3.0$ (shown in red) and $\beta_{\text{corr}} = \hat{\beta} = +10.0$ (shown in black) resulting in positive correlations.

galaxy in each $\hat{\beta}$ corrected data-set and used to estimate the random variable, ζ . It is obvious that for $\beta_{\text{corr}} = -10.0$, the curving of the M - Z distribution produces a sizeable gap between the majority (distant) of galaxies and the corresponding $m_* = 19.8$ mag line which therefore manifests as the anti-correlation observed in left panel of Figure 9.9. However, for $\beta_{\text{corr}} = +10.0$ the M - Z distribution is now outwardly curved resulting in the faintest galaxies being the most distant ones giving an $m_* = 21.9$ mag (shown in green). Therefore, when we estimate ζ in this case, the curving of the distribution in this way (top-right region) dominates the ζ calculation (since the majority of galaxies are located here). We can see that these galaxies are not distributed as far from the m_* limit as with the cases of $\beta_{\text{corr}} = +10.0$ does and therefore the resulting correlation in the right-hand panel of Figure 9.9 is diminished.

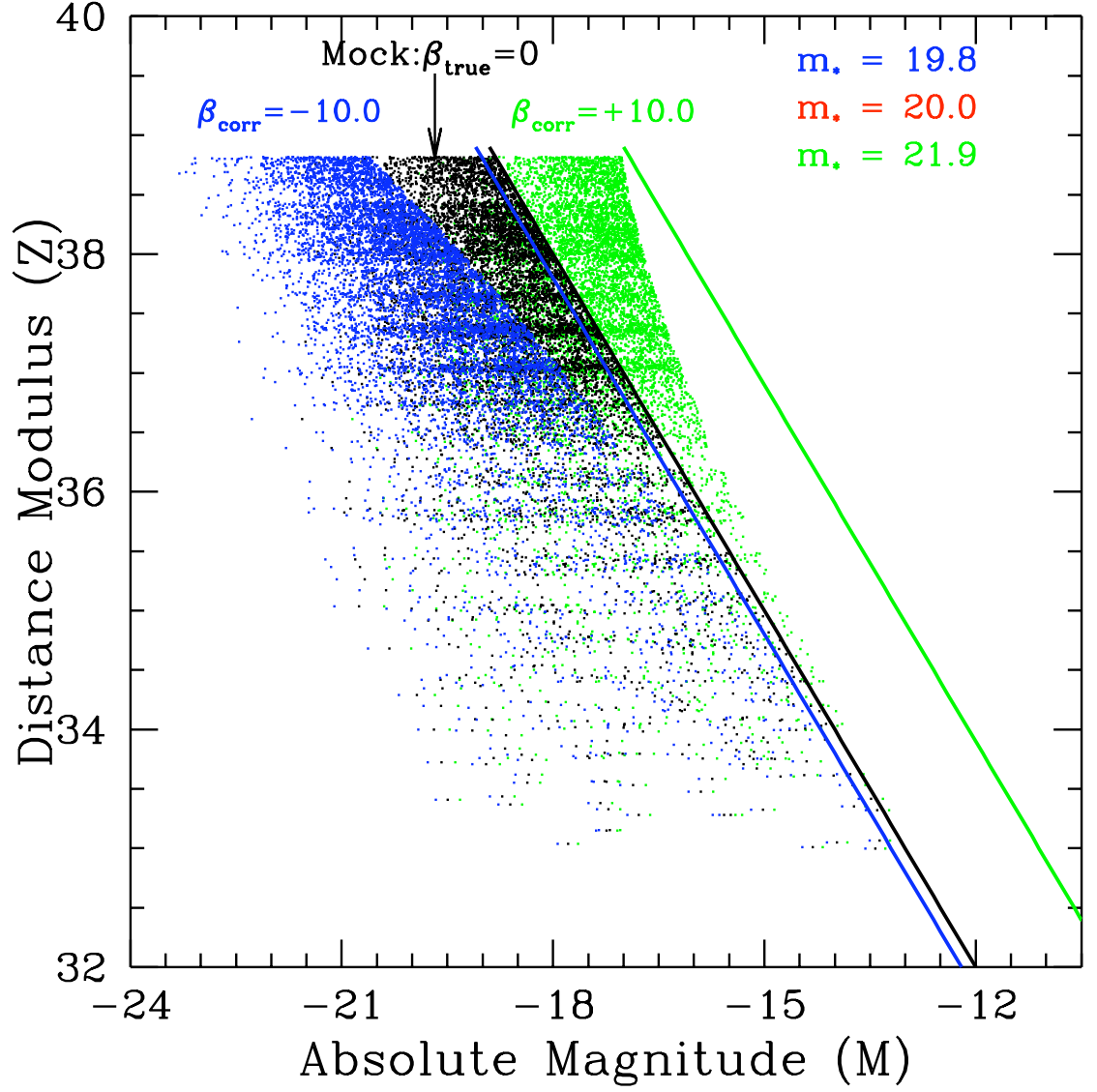


Figure 9.10: M - Z distribution for the MGC mock catalogue examined in Figure 9.9 with $\beta_{\text{true}} = 0$ (shown as the black distribution of points). The other coloured distributions illustrate the effect of correcting the absolute magnitudes according to our trial $\hat{\beta}$ correction. In this example the blue distribution represents a correction of $\beta_{\text{corr}} = \hat{\beta} = -10.0$ with a resulting $m_* = 19.8$ mag. The green points represent $\beta_{\text{corr}} = \hat{\beta} = +10.0$ with a resulting $m_* = 21.9$ mag.

9.2.3 SDSS-Mocks

We now use the JTH method for estimating ζ and χ and use Figure 9.11 as a guide for choosing the optimum value of δZ where we have plotted $(\hat{\beta}, \rho)$ curves for $\delta Z = 0.01, 0.1, 0.5$ and 1.0 . We can see in the top two panels the same characteristic noise dominated behaviour already observed in the MGC analysis for small value of δZ . Although a $\delta Z = 0.5$ would appear an adequate choice we do observe odd behaviour in ρ across the range in $\hat{\beta}$. In the bottom panels of Figure 9.11 it is clear to see that ρ peaks sharply at $\hat{\beta} = 0.0$ for all the mock realisations. However, the value of ρ at this point is ~ 0.07 at $\delta Z = 0.5$ and ~ 0.157 at $\delta Z = 1.0$. Furthermore, there appears to be two possible range in $\hat{\beta}$ where $\rho = 0$ for which ζ and χ would be considered to be unocorrelated. For $\delta Z = 0.5$ this occurs at $-0.5 \lesssim \hat{\beta} \lesssim -0.4$ and $0.2 \lesssim \hat{\beta} \lesssim 0.8$. For $\delta Z = 1.0$ there is a slightly broader spread at $-1.7 \lesssim \hat{\beta} \lesssim -1.1$ and $1.1 \lesssim \hat{\beta} \lesssim 1.8$.

When we apply ρ to evolved mocks in Figure 9.12, where $\beta_{true} = 1.0, 2.0$ and 3.0 , we observed exactly the same behaviour. For each $\beta_{true} = 1.0, 2.0$ and 3.0 , ρ shows anti-correlation from most large negative values, then crosses through $\rho = 0$ and peaks at the expecting β correction value as indicated by the vertical red lines in each case. The reason for this is not immediately clear, however, in the following sub section we offer possible solutions that indicate the ρ estimator is perhaps insufficient for this type of analysis.

9.2.3.1 Errorr Analysis

As with MGC in § 9.2.2.1 we use 200 bootstraps for the SDSS analysis and interpolate each $\hat{\beta}$ curve within a small range in intervals of $\hat{\beta} = 0.05$. As in Figure 9.12 we have applied a δZ width of 0.5 . In Figure 9.13 we have plotted the $(\hat{\beta}, \rho)$ curves for $\beta_{true} = 0.0$ and 1.0 . As we have already observed in Figure 9.12, for each β_{true} there are *two* distinct points where $\hat{\beta}$ crosses $\rho = 0$. If we look at Figure 9.14 we can observe this distribution for each case. For $\beta_{true} = 0$ in the top panel of Figure 9.14 we can see that the distribution left of $\hat{\beta} = 0$ is tightly distributed with a mean value of $\bar{\beta} = 0.41$ and a corresponding 68% CI of $[0.382, 0.451]$. For the distribution right of $\hat{\beta} = 0$, we observe a $\bar{\beta} = -0.53$ and a CI = $[-0.461, -0.600]$. In this case the overall distribution is broader than its left counterpart. In the bottom panel of Figure 9.14 we see an almost identical distribution for $\beta_{true} = 1.0$. In this case the $\bar{\beta}(\rho = 0) = 1.416$ and 0.472 with respective CI's = $[1.382, 1.451]$ and $[0.399, 0.538]$.

In Figures 9.15 and 9.16 we show similar distributions for $\beta_{true} = 2.0$ and 3.0 . Once again we see, particularly in Figure 9.16, very similar distributions for $\hat{\beta}(\rho = 0)$ as in

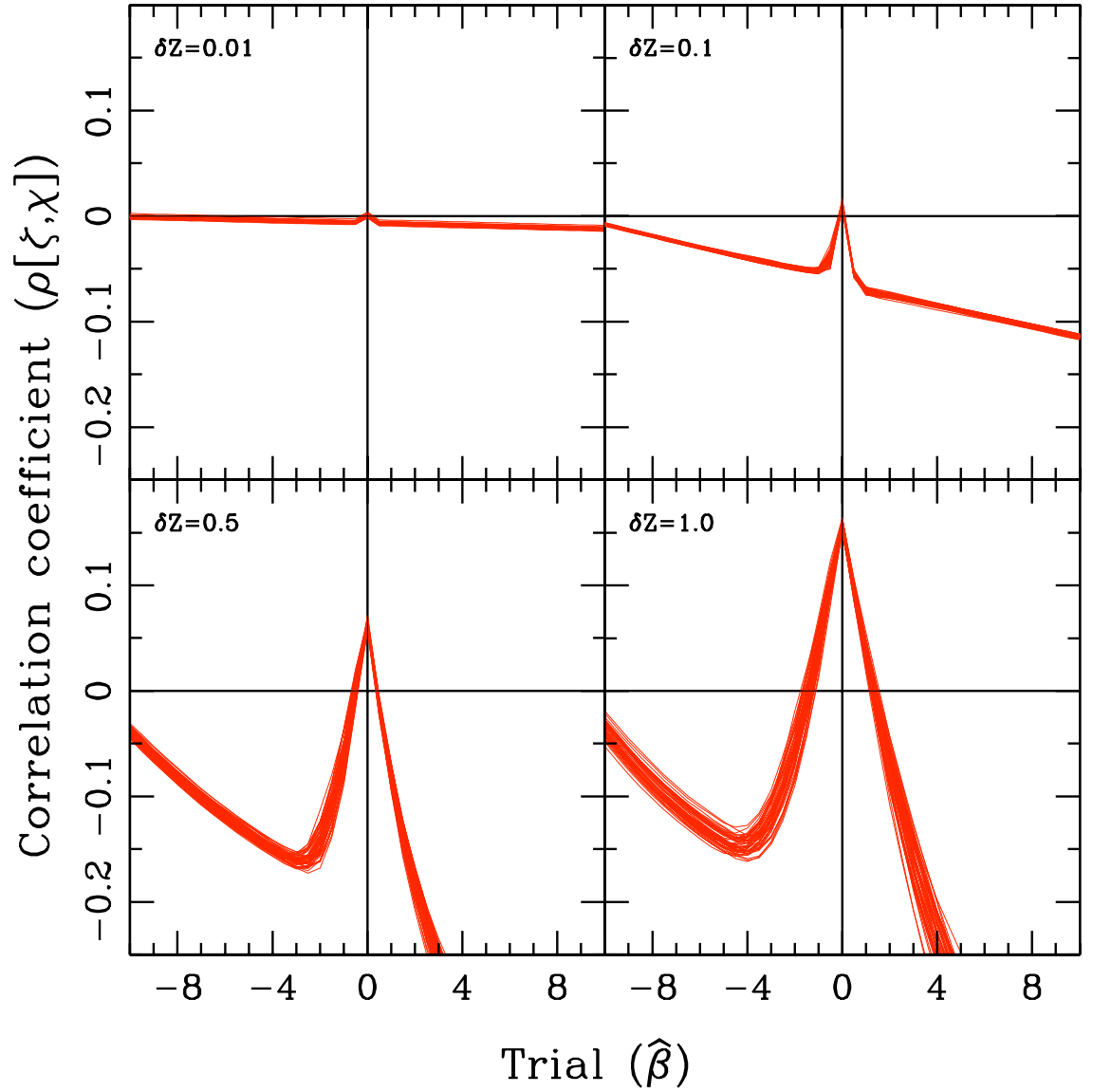


Figure 9.11: Four plots demonstrating the shot noise behaviour in the SDSS mocks for varying sizes of δZ . In this example we applied $\delta Z = 0.01, 1.0, 0.5$ and 1.0 . As with MGC in Figure 9.5 lower values of $\delta Z = 0.01$ and 0.1 appear to be insufficient for ρ to constrain β . However, although $\delta Z = 0.5$ (bottom-left) and indeed 1.0 (bottom-right) now appear to be acceptable choices we observed behaviour in ρ that is inconsistent with the true value of β for the evolved mock catalogues.

Figure 9.14. For $\beta_{true} = 2.0$ we have determined $\bar{\beta} = 2.416$ and 1.472 with a respective 68% CI=[2.381, 2.499] and [1.399, 1.538]. For $\beta_{true} = 3.0$ we have determined a corresponding $\bar{\beta} = 3.416$ and 2.472 with a respective 68% CI=[3.381, 3.449] and [2.399, 2.538].

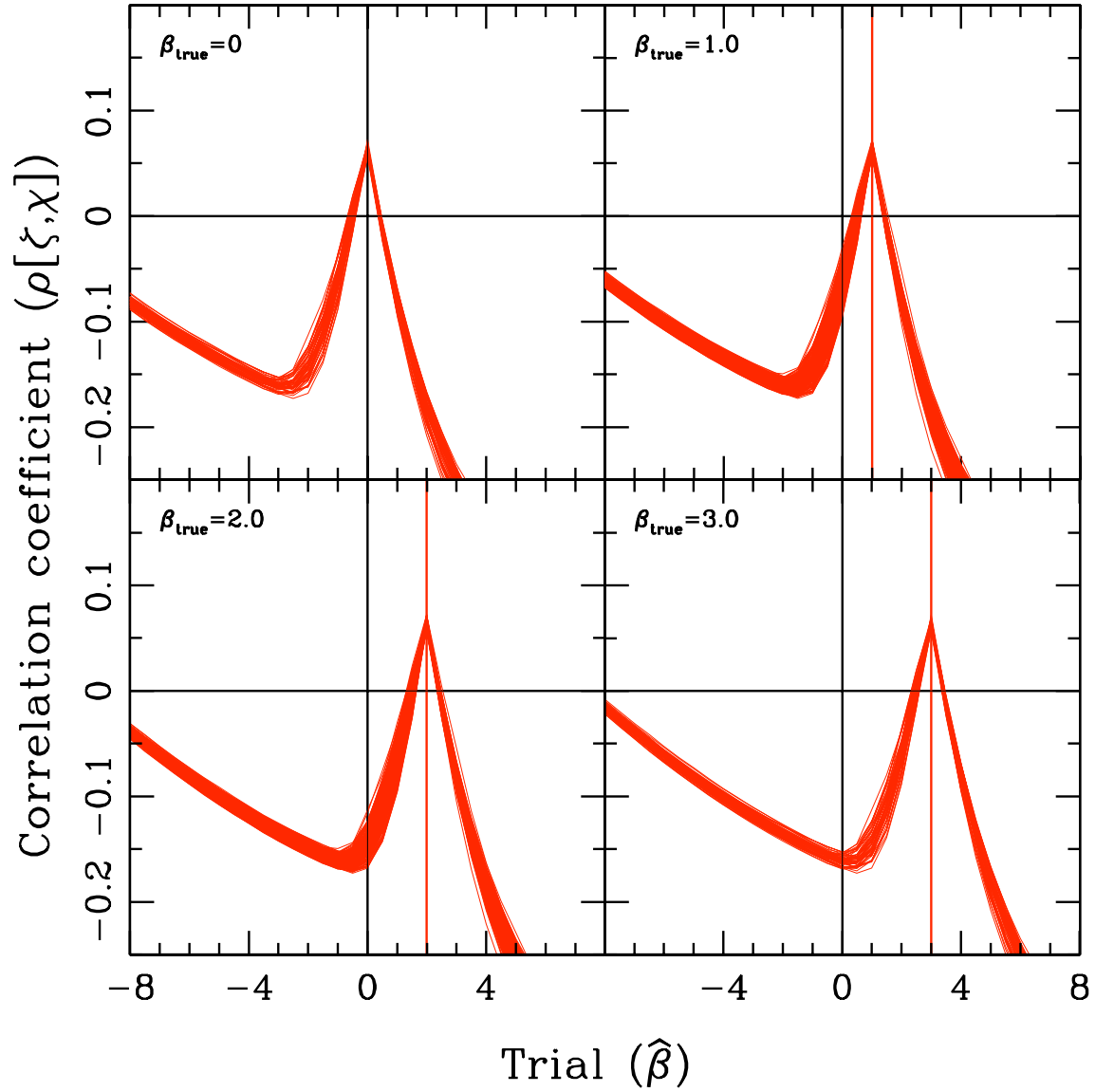


Figure 9.12: Plots showing the resulting correlation coefficient, ρ , for SDSS mock catalogues generated from an evolving LF with $\beta_{\text{true}}=0.0, 1.0, 2.0, 3.0$. In all four cases we have fixed $\delta Z = 0.5$.

So we are left with the question of why do we see, in each β_{true} , two distinct ranges for constraining $\hat{\beta}$? Whilst there has been a rigorous effort to check for inconsistencies in the programming, none as yet has been found. It may perhaps be necessary to explore the higher order moments of ρ in order to find some answers. However, at this juncture we shall apply the entropy approach that has all the higher order moments contained within it. As we shall see this has proven to be crucial for constraining evolution.

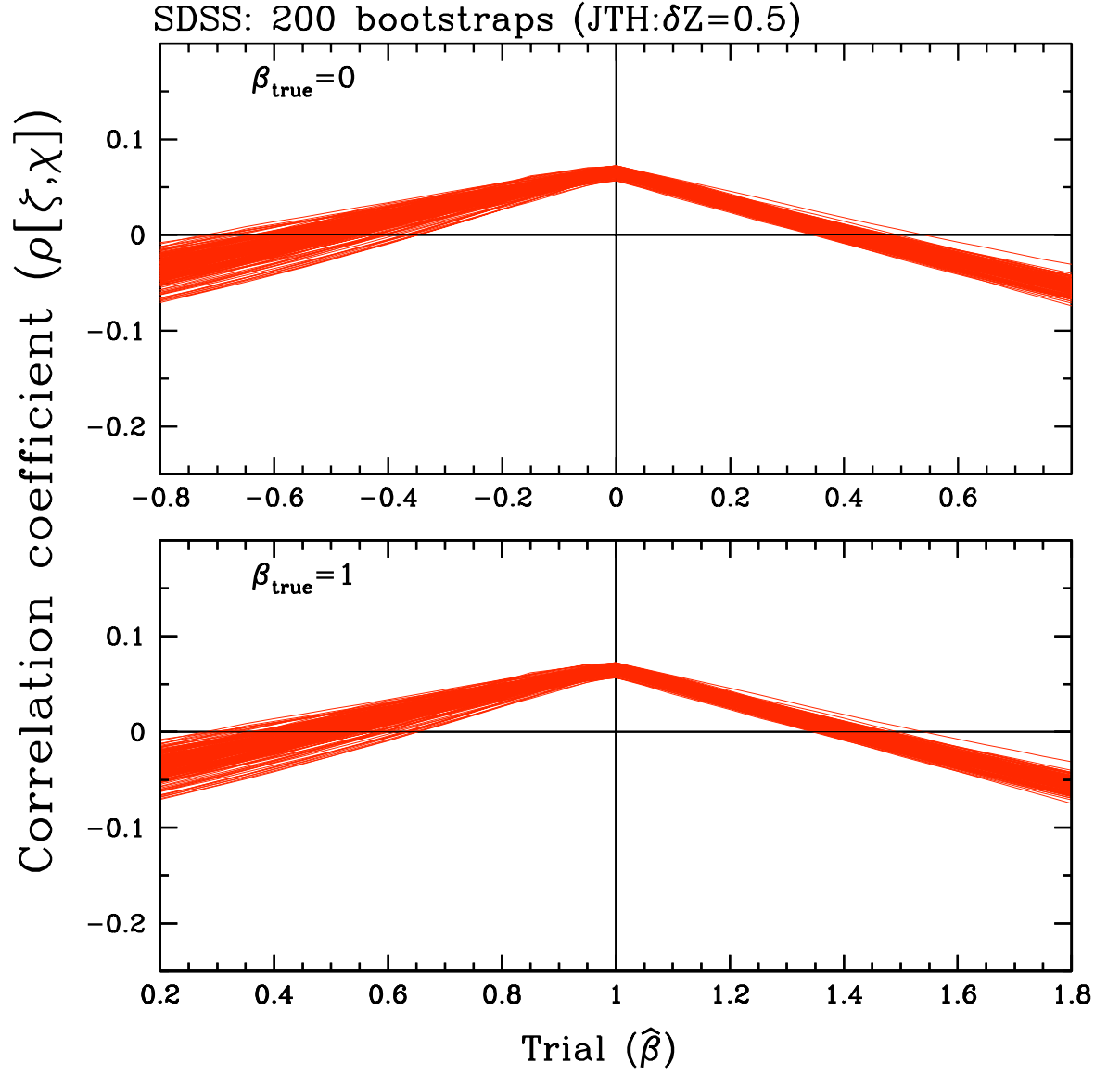


Figure 9.13: ρ versus $\hat{\beta}$ for 200 SDSS bootstraps with $\beta_{true}=1.0$ and 2.0 . As already observed in Figure 9.12 we can see that $\hat{\beta}$ crosses $\rho = 0$ at two distinct points either side of β_{true} . Figure 9.14 shows the corresponding $\hat{\beta}(\rho = 0)$ distribution for these two panels.

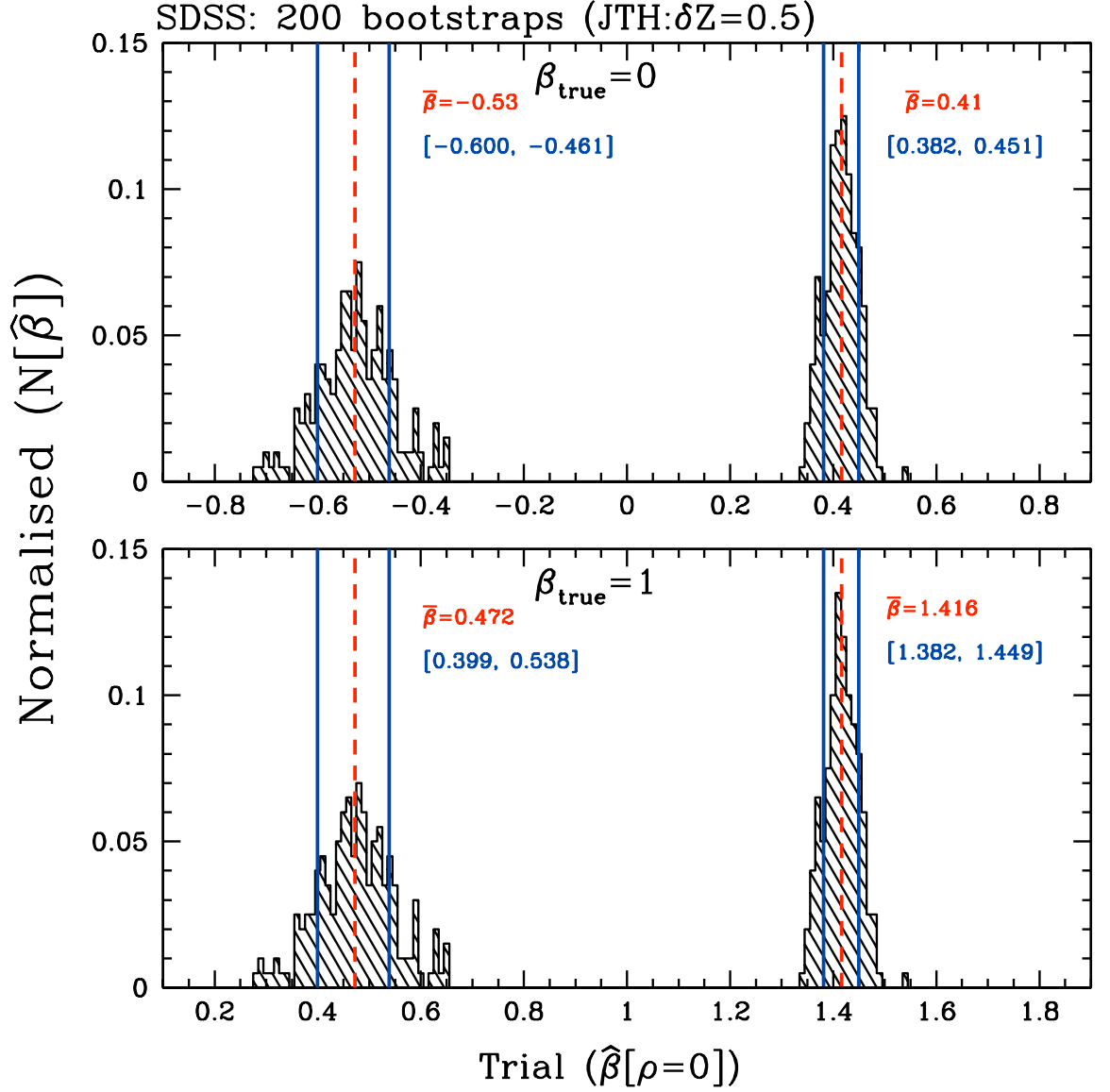


Figure 9.14: $\hat{\beta}(\rho = 0)$ distribution for 200 SDSS bootstraps at $\beta_{\text{true}} = 0.0$ and 1.0 . In each panel the red dashed line represents the mean of the distribution, with the corresponding value, $\bar{\beta}$, in red. The blue vertical lines show the 68% confidence interval with the corresponding values in blue parenthesis. Unlike MGC, the SDSS mocks indicate two possible values for β_{true} which for $\beta_{\text{true}} = 0$ is $\bar{\beta} = 0.41$ and -0.53 , and for $\beta_{\text{true}} = 1$ is $\bar{\beta} = 1.416$ and 0.472 .

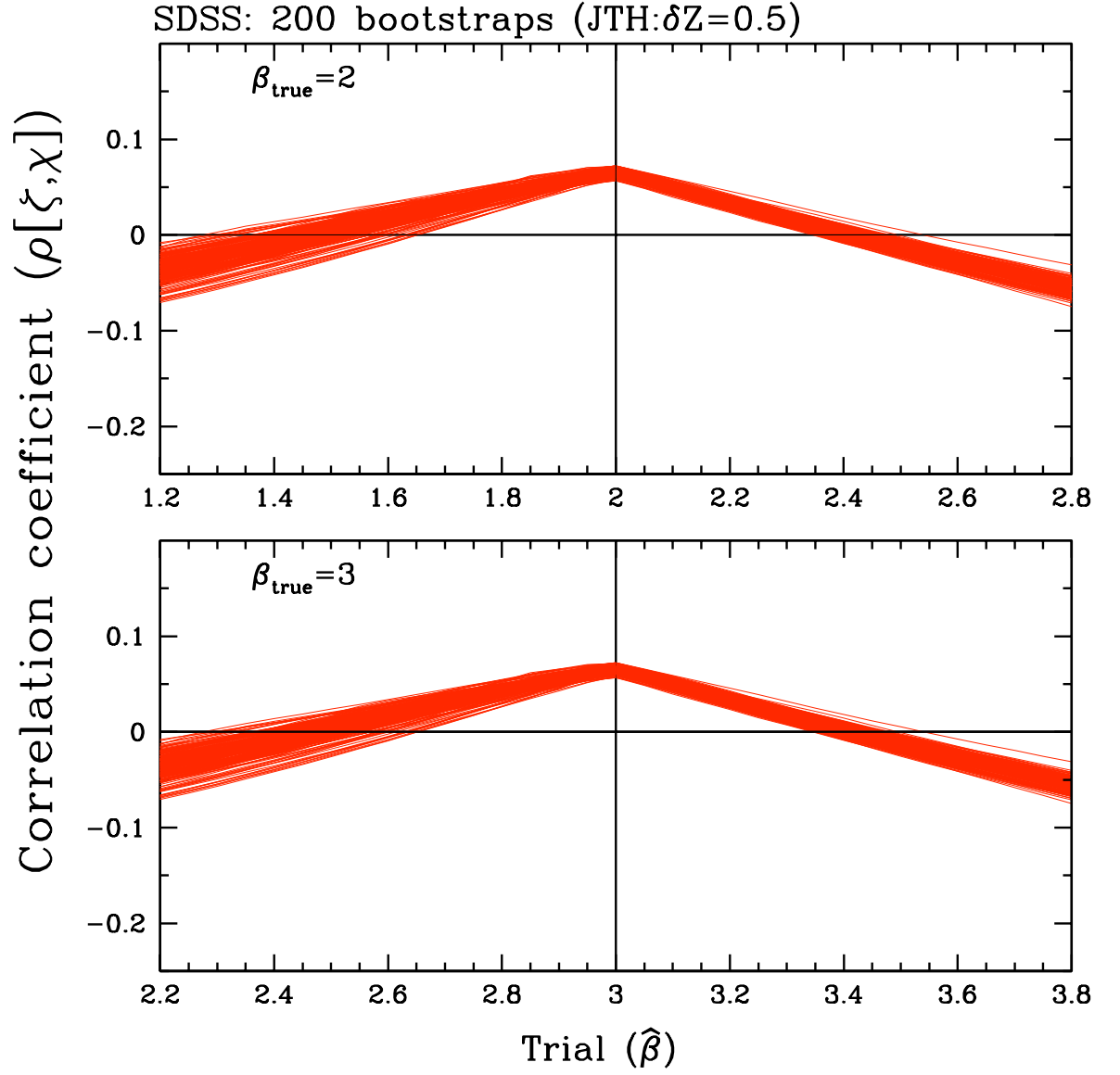


Figure 9.15: ρ versus $\hat{\beta}$ for 200 SDSS bootstraps with $\beta_{true}=2.0$ and 3.0 . As already observed in Figure 9.12 we can see that $\hat{\beta}$ crosses $\rho = 0$ at two distinct points either side of β_{true} . Figure 9.14 shows the corresponding $\hat{\beta}(\rho = 0)$ distribution for these two panels.

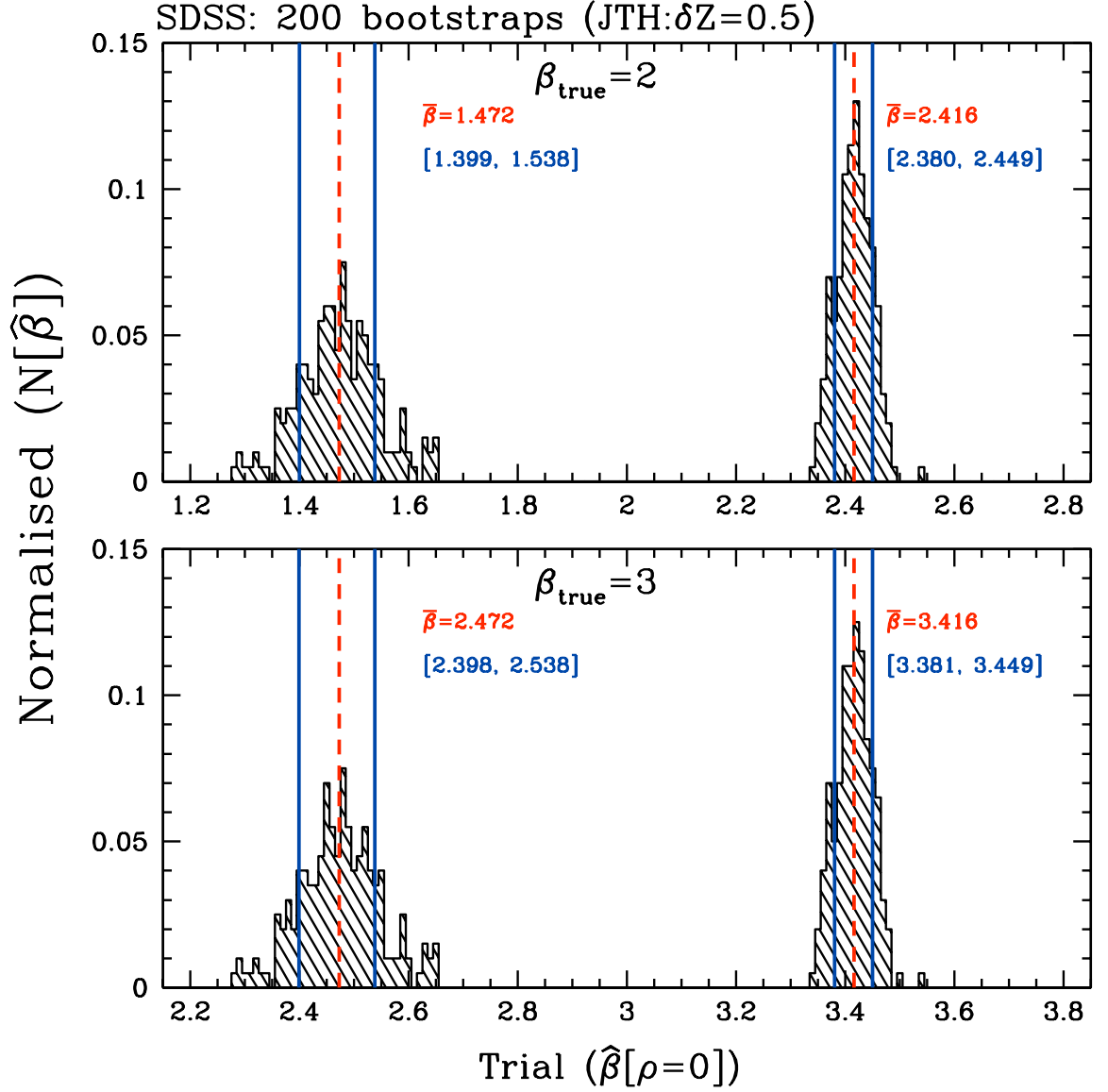


Figure 9.16: $\hat{\beta}(\rho=0)$ distribution for 200 SDSS bootstraps at $\beta_{\text{true}} = 2.0$ and 3.0 . In each panel the red dashed line represents the mean of the distribution, with the corresponding value, $\bar{\beta}$, in red. The blue vertical lines show the 68% confidence interval with the corresponding values in blue parenthesis. Unlike MGC, the SDSS mocks indicate two possible values for β_{true} which for $\beta_{\text{true}} = 2$ is $\bar{\beta} = 2.416$ and 1.472 , and for $\beta_{\text{true}} = 3$ is $\bar{\beta} = 3.416$ and 2.472 .

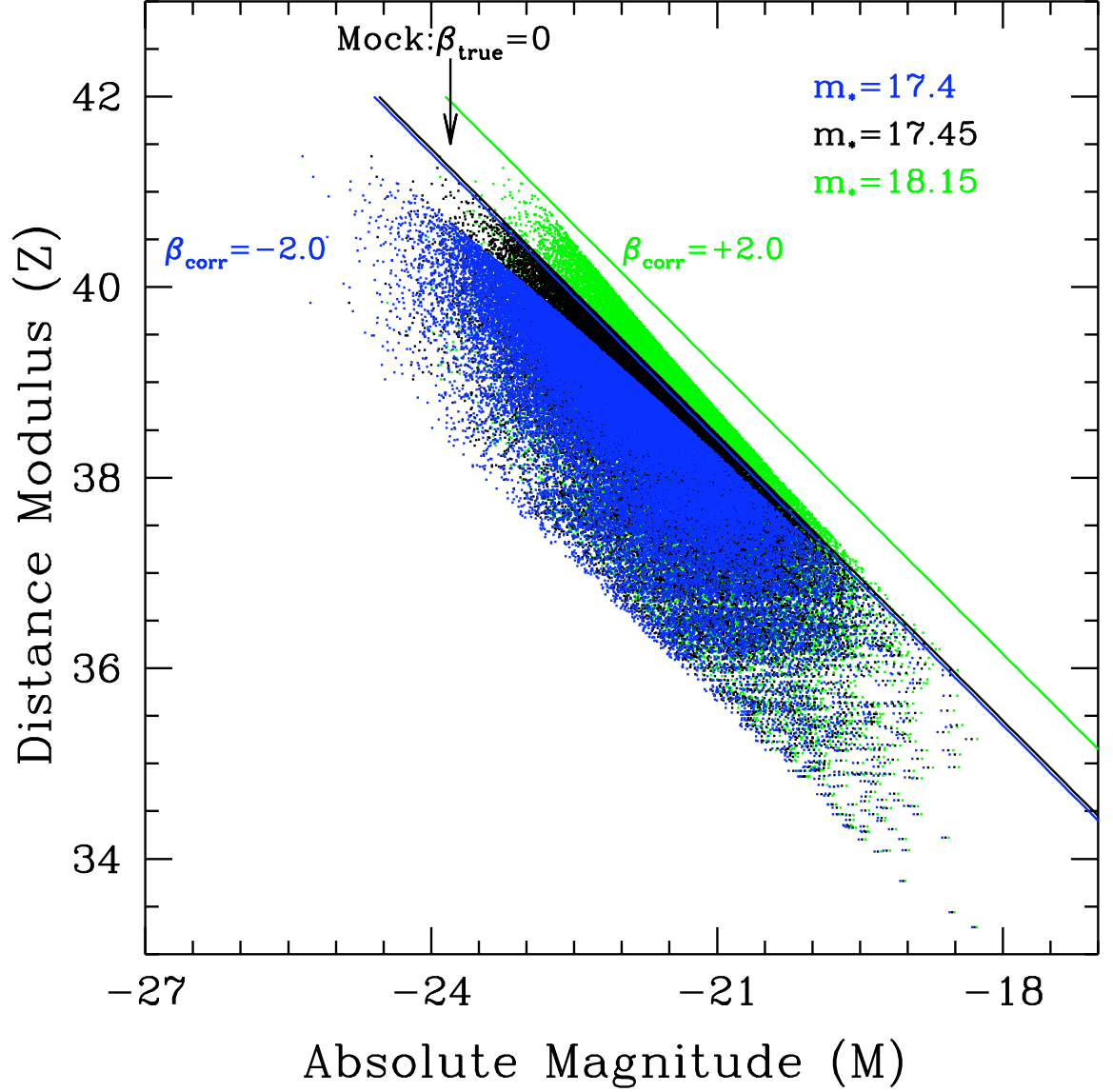


Figure 9.17: M - Z distribution for SDSS mock with $\beta_{\text{true}} = 0$ shown as the black distribution of points. The other coloured distributions illustrate the effect of correcting the absolute magnitudes according to our trial $\hat{\beta}$ correction. In this example the blue distribution represents a correction of $\beta_{\text{corr}} = \hat{\beta} = -2.0$ with a resulting $m_* = 17.4$ mag. The green points represent $\beta_{\text{corr}} = \hat{\beta} = +2.0$ with a resulting $m_* = 18.15$ mag.

9.3 Results Part II - Relative Entropy

9.3.1 MGC-Mocks (R01)

Before we test mocks that are drawn from an evolving LF, let us firstly consider the case when $\beta_{true} = 0$ resulting in a mock drawn from a Universal LF. In Figure 9.18 we use these mocks to observe the behaviour of the relative entropy for varying resolutions of the imposed mesh on the ζ - χ plane. We try a range $-10.0 \leq \hat{\beta} \leq 10.0$ in increments of 0.5. For each cell size we use 10 mock realisations. For each trial cell size we expect the relative entropy to minimise at a $\hat{\beta}_k = \beta_{true}$. As we can see in Figure 9.18 this is indeed the case! However we make the following observations regarding the choice in cell size:

1. For the cell sizes where we have a total of 25 cells (for a 0.2×0.2 cell), and 100 cells (for the 0.1×0.1 cell), the relative entropy does minimise, but for a possible range of $\hat{\beta}$. For the 0.2×0.2 cell size this range is $1.4 \lesssim \hat{\beta} \lesssim 0.8$, and for the 0.1×0.1 cell size we observe a narrower range of $-0.6 \lesssim \hat{\beta} \lesssim 0.2$. As we move to a smaller cell size of 0.02×0.02 the relative entropy appears to minimise very sharply at the exact true value for β . For even smaller cell sizes at 0.001×0.001 (corresponding to 1,000,000 cells on a unit square) we observe that although there is a sharp minimisation at $\hat{\beta} = 0$, the overall entropy curve across the interval $[-10, 10]$ is flattened.
2. Recalling that for the condition where the observed, p_i , equals the theoretical model, q_i , the relative entropy, S will equal zero. In Figure 9.18, as the number cells increases, the observed relative entropy moves further from $S = 0$.

In the case of our first observation, if we have a small number of cells making up our mesh (meaning larger and larger individual cell sizes), it becomes apparent that the range for which we can minimise S broadens. Conversely if we have a grid size that has a very large number of cells then we would expect the result to be noise dominated since there will only be a few galaxies in every cell. This would explain why we see S begin to flatten for cell size of 0.001×0.001 but shows very little variation for a broad range of $\hat{\beta}$ for 0.2×0.2 cells (Figure 9.18).

The second observation relates to the theoretical model, q_i , that we have chosen to use. For our model we have simply distributed the points on the ζ - χ uniformly such that there is an equal number in each cell. Although this represents an ideal case it does not of course necessarily represent a realistic model since it does not account for

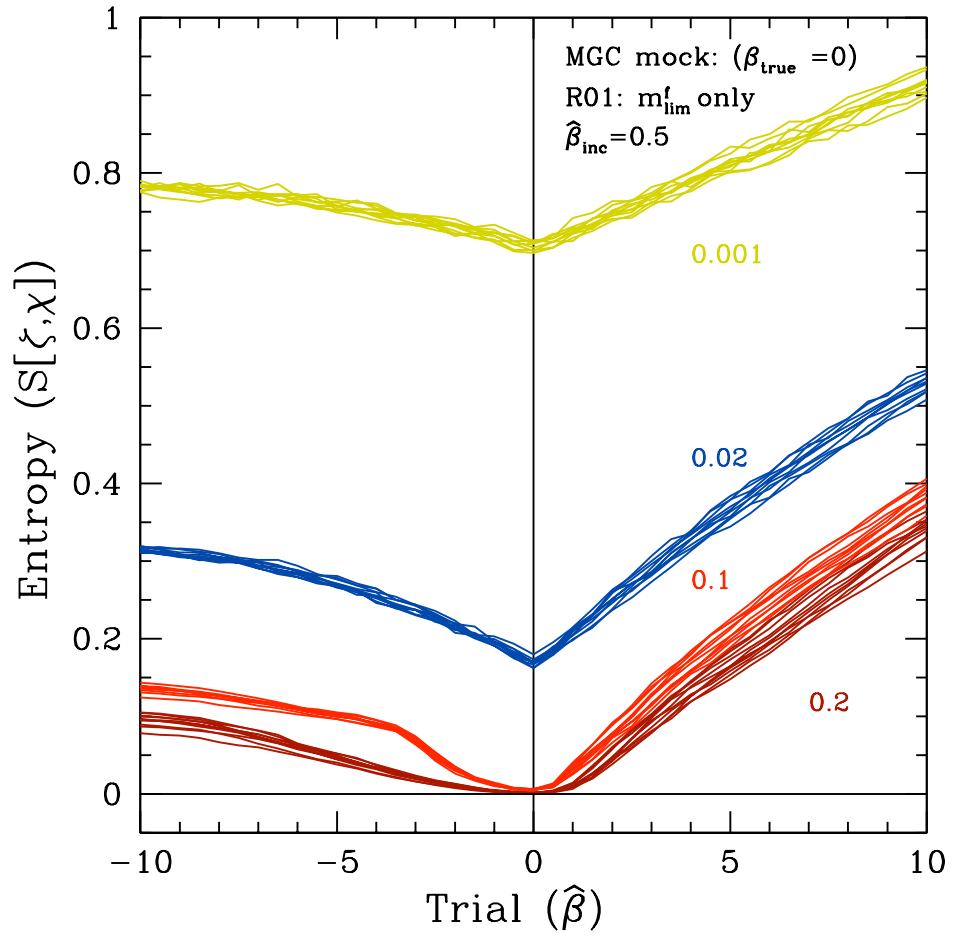


Figure 9.18: Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks for varying mesh sizes applying the R01 method. In this example we have used 10 mock realisations for each mesh size drawn from a Universal LF i.e $\beta_{true} = 0$. For each mesh the cell sizes are shown as follows: 0.001×0.001 in green, 0.02×0.02 in blue, 0.1×0.1 in red and 0.2×0.2 in magenta. The main result from this plot indicates that despite varying resolutions in the mesh, the relative entropy minimises at the correct trial $\hat{\beta} = \beta_{true} = 0$.

any sample variance that would result from a real galaxy distribution. Therefore, on small scales, the resulting ζ - χ distribution will inevitably exhibit clustering which will alter the relative entropy compared to our theoretical model. Hence, in Figure 9.18 we observe a shift in S from zero on small grid scales. Although we could use a more realistic theoretical model, it is important to note that the model we adopt does not seem adversely affect its ability to constrain β_{true} .

Continuing for now with the R01 method estimating ζ , we now apply our mocks that are drawn from an evolving LF. We apply successive values of $\beta_{true} = 1.0, 2.0$ and 3.0 and assume a mesh size cell size of 0.02×0.02 as show in Figure 9.19. For each β_{true}

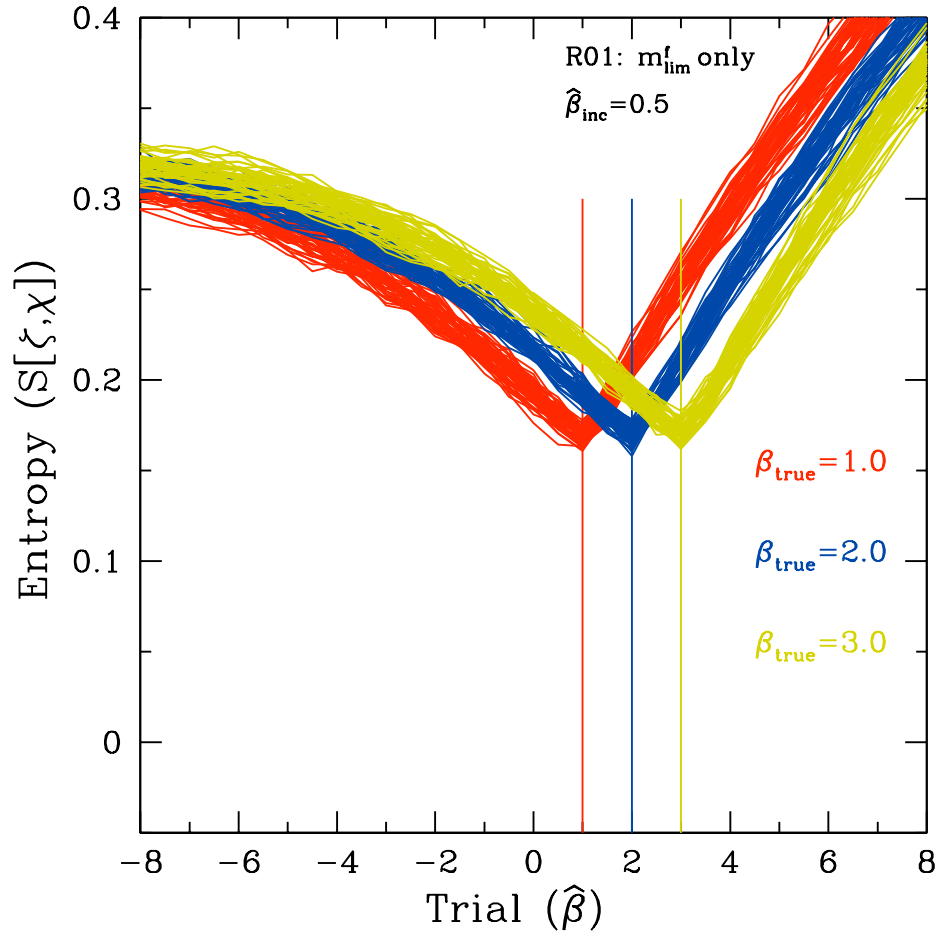


Figure 9.19: Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks for three different values of β_{true} applying the R01 method. For each trial $\hat{\beta}$ we have applied 50 mock realisations with successive values of $\beta_{\text{true}} = 1.0, 2.0$ and 3.0 . In each case we can clearly see that the relative entropy approach minimises at the correct value. Notice also that values chosen by S are the negative of the β_{true} value. This is expected since we correcting the data in such a way as to render the M - Z distribution in a separable form.

we have generated 50 mock realisations. The resulting entropy, S , show that for mocks generated with evolution of the order $\beta_{\text{true}} = 1.0, 2.0$ and 3.0 , $S(\zeta, \chi)$ minimises at respective $\hat{\beta} = -1.0, -2.0$ and -3.0 .

9.3.2 MGC-Mocks (JTH)

As we have already demonstrated in Chapters 4 and 6, one of the features of the JTH generalisation allows us to choose a suitable width, δZ , for the S_1 and S_2 regions in order that the regions remain separable within the bright and faint apparent magnitude limits of a given survey. Computationally, this implies that the smaller the width of δZ ,

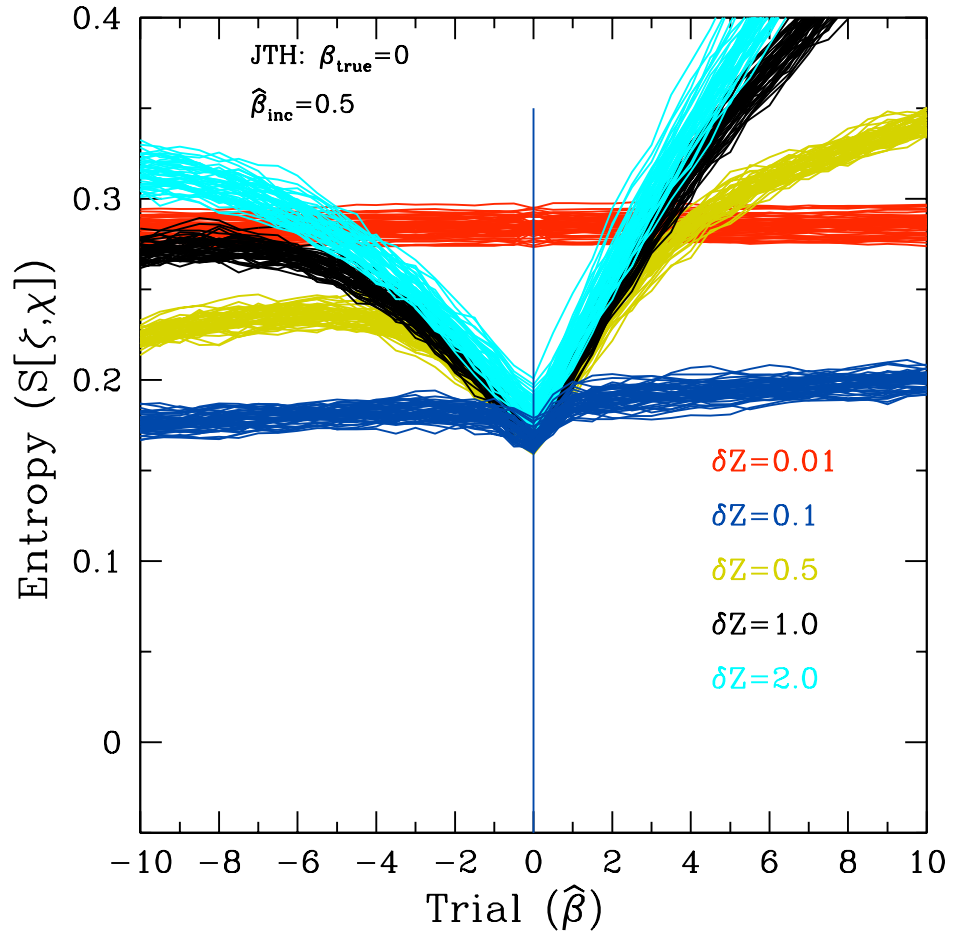


Figure 9.20: Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks implementing the JTH for different values of δZ . Similarly with the correlation coefficient, we observe that for a $\delta Z = 0.01$, the entropy S does not minimise. However, for values ≥ 0.5 we can observe a clear signal in S . It is also worth noting that at $\delta Z = 2.0$, the value at which S minimises is unchanged unlike in the same case for the correlation coefficient in Figure 9.4. Although the number of galaxies going into the ζ and χ estimation is changing with the size of δZ , the number of galaxies going into the theoretical model changes with it so that we are comparing p and q with an equal total number, N .

the faster the code will run for each estimation of the random variable, ζ . However, as revealed in Chapter 6, as we move to increasingly smaller values of δZ , the shot noise begins to dominate and our ability to recover any meaningful signal diminishes as a result. Therefore, as we incorporate the JTH method in to our analysis in this section, we experiment with a range of δZ for the ζ estimation that will maximise both the processing time and amount of signal required to constrain $\hat{\beta}$. In exactly the same way as in Chapter 4 when applying the JTH method we assume a bright limit equal to the apparent magnitude of the brightest galaxy for every $\hat{\beta}$ corrected sample. Figure 9.20

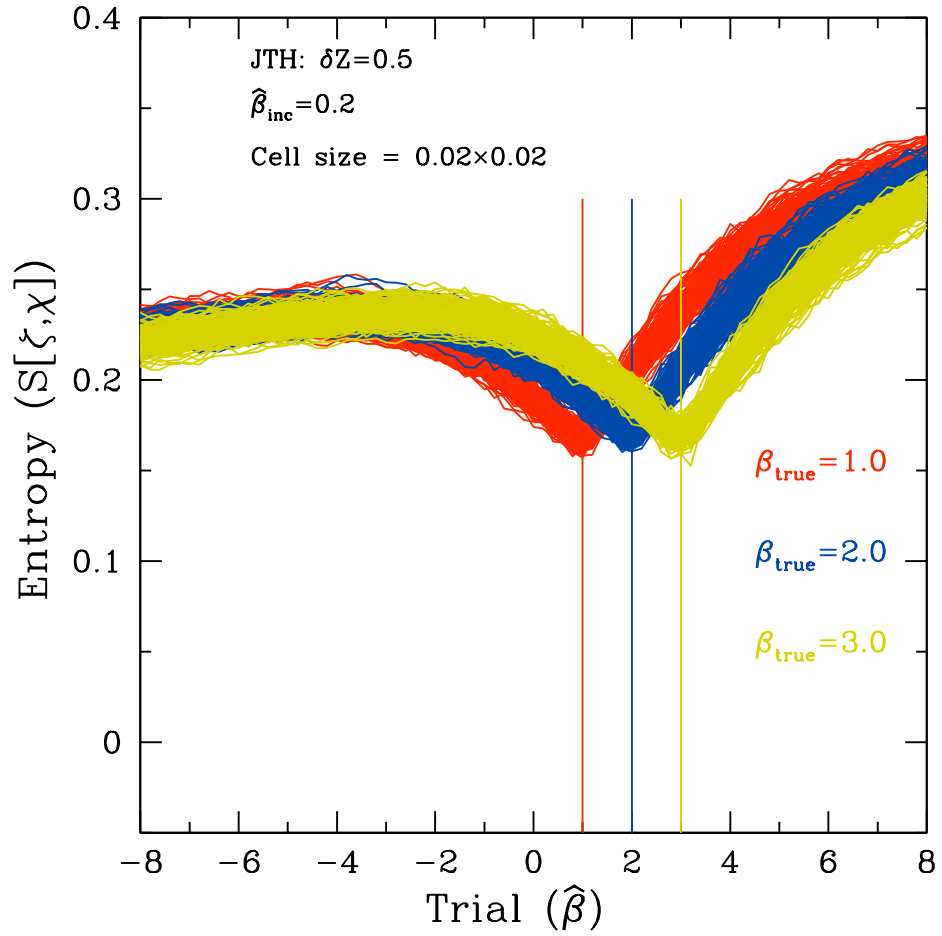


Figure 9.21: Relative entropy, S versus trial $\hat{\beta}$ using MGC mocks for $\beta_{true} = 1.0, 2.0$ and 3.0 applying the JTH method with a $\delta Z = 0.5$. We have imposed a mesh cell size of 0.02×0.02 and have incremented trial $\hat{\beta}_{inc} = 0.2$. Just as with the R01 method we observe that in each case, applying the JTH method, the relative entropy estimator appears to constrain β extremely well.

illustrates the effect of varying δZ for values, 0.01 (red), 0.1 (blue), 0.5 (green), 1.0 (black) and 2.0 (dark green). We have assumed in each case a mock with $\beta_{true} = 0$ and generate 50 mock realisations in each case. For a $\delta Z = 0.01$ we observe a near flat-line for the full range of $\hat{\beta}$ with only a minor kink as $\hat{\beta}_k$ pass zero. Although as δZ increases to 0.1 we can observe the relative entropy minimise at $\hat{\beta} = 0$, the signal remains quite weak. It's not until $\delta Z \gtrsim 0.5$ that we achieve a signal which is similar to that shown in Figure 9.19 for the R01 method.

In Figure 9.21 we have generated 200 mock realisations for $\beta_{true} = 1.0, 2.0$ and 3.0 and with apply a $\delta Z = 0.5$ throughout. Once again we can see that the relative entropy's in each case minimise exactly at the respective expectant values of $\hat{\beta} = -1.0$,

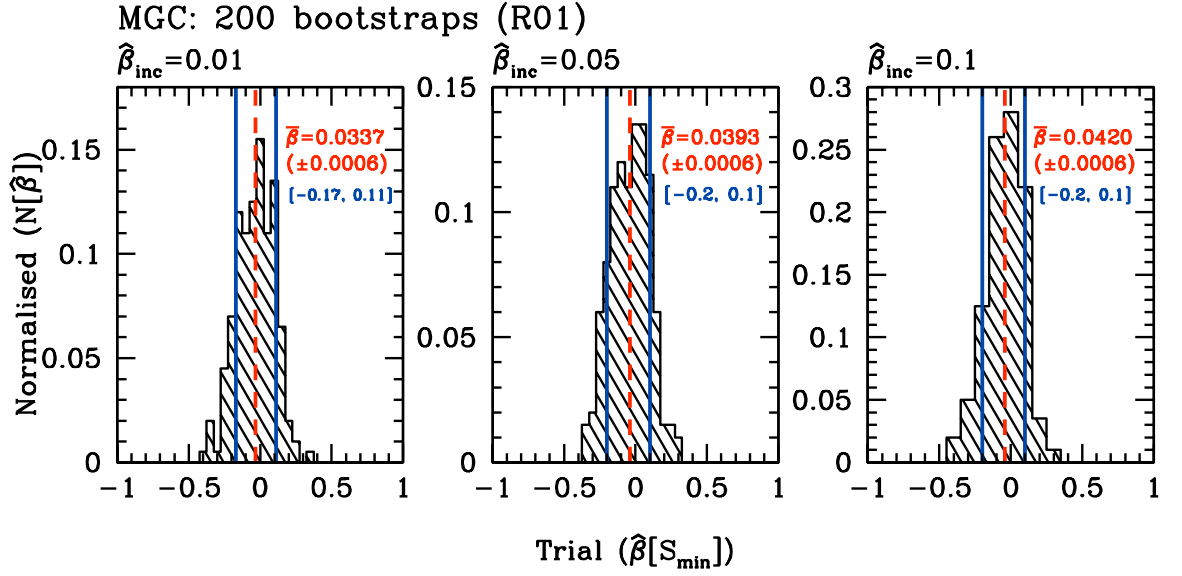


Figure 9.22: $\hat{\beta}(S_{\min})$ distribution for varying increments of $\hat{\beta}$. In this figure we want to determine a suitable increment in $\hat{\beta}$ for finding the minimised entropy, S_{\min} . We have tested values $\hat{\beta}_{\text{inc}} = 0.01, 0.05$ and 0.1 . The mean of each distribution, $\bar{\beta}$, is shown as a dashed red line with the 68% CI indicated by the parallel blue lines in each panel. Since there is a difference of only 0.005 in $\bar{\beta}$ between $\hat{\beta}_{\text{inc}} = 0.01$ and 0.05 we have chosen to use increments of 0.05 to maximise computational time.

-2.0 and -3.0.

9.3.2.1 Error Analysis

To assess the error distribution of $\hat{\beta}$ we applied the same number of bootstrap realisations as with the ρ statistic. To further optimise the runtime of the code we initially tested 3 values of $\hat{\beta}_{\text{inc}}$ to observe whether there was any significant differences between three distributions. The increments we chosen were $\hat{\beta}_{\text{inc}}=0.01, 0.05$ and 0.1 and the corresponding distributions are shown in Figure 9.22. Once again we have superimposed the 68% CI in blue and the mean value of $\hat{\beta}$ as the red dashed line. In this example $\beta_{\text{true}}=0$. As we can see all three incremental values are all very tightly distributed about zero with only variation of 0.0083 in the mean from $\hat{\beta}_{\text{inc}}=0.01$ to 0.1 and only 0.005 difference between 0.01 and 0.05 . Therefore, we chose to increment $\hat{\beta}_{\text{inc}}$ in intervals of 0.05 .

Figure 9.26 shows the resulting $\hat{\beta}(S_{\min})$ distributions for the trial $\beta_{\text{true}}=0, 1, 2$ and 3 . We can see that in all cases the distributions are very narrow compared to the ρ test, with very little visual evidence of skewness. For $\beta_{\text{true}}=0$ we recover a mean value

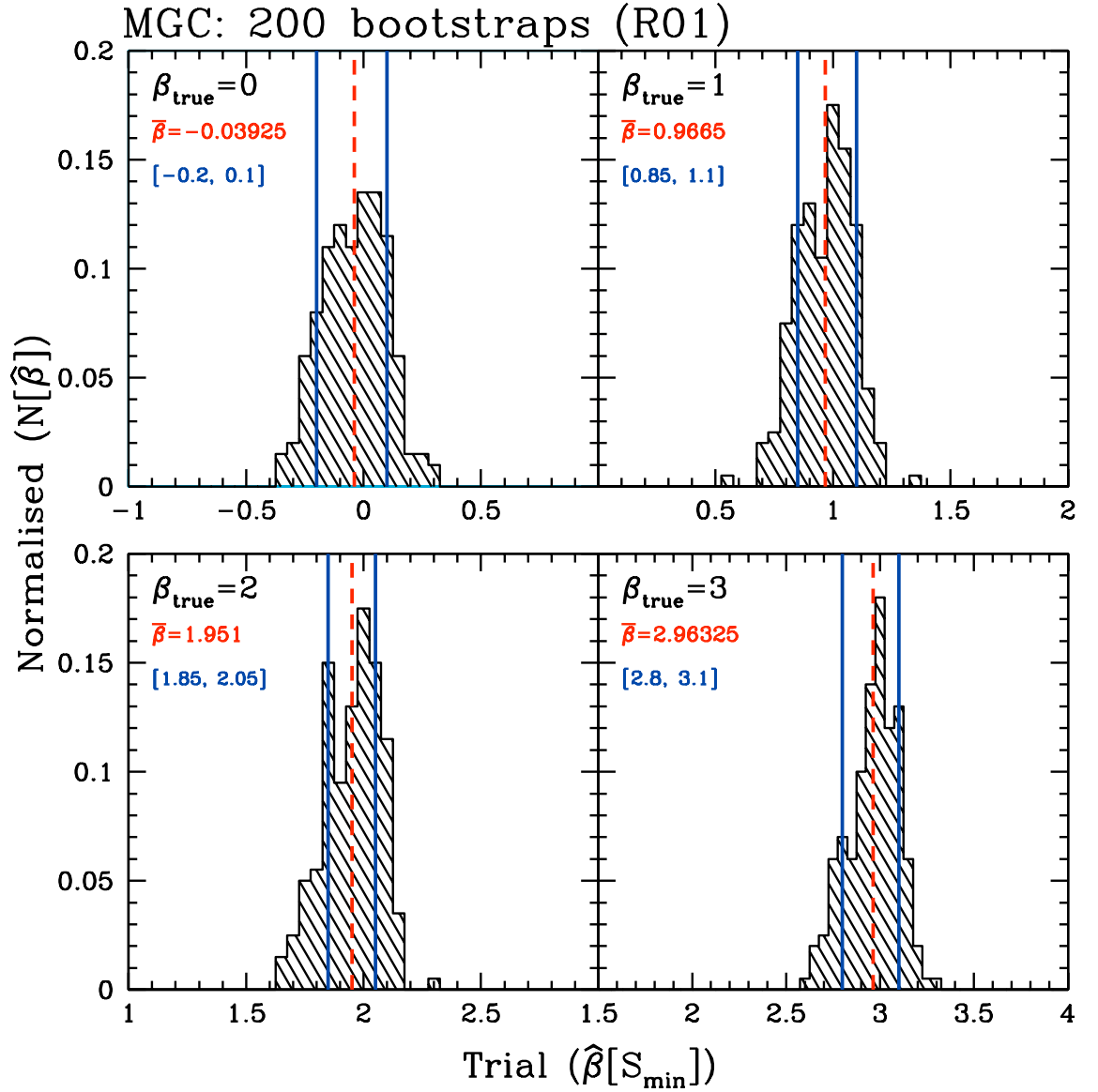


Figure 9.23: $\hat{\beta}(S_{\min})$ distribution for 200 MGC bootstraps for $\beta_{\text{true}} = 0, 1.0, 2.0$ and 3.0 . We follow the same convention as with the ρ testing superimposing the 68% CI's solid blue vertical lines, and the average of the distribution as a vertical red dashed line.

of $\bar{\beta} = -0.03925$ with a 68% CI of $[-0.2, 0.1]$; $\beta_{\text{true}} = 1$, $\bar{\beta} = 0.9665$ with a 68% CI of $[0.85, 1.1]$; $\beta_{\text{true}} = 2$, $\bar{\beta} = 1.951$ with a 68% CI of $[1.85, 2.05]$; and $\beta_{\text{true}} = 3$, $\bar{\beta} = 2.96325$ with a 68% CI of $[2.8, 3.1]$.

As we have discussed already in this chapter, the spread of these distributions would vary if one were to vary the cell widths for the entropy calculation. However, we have found that the cell size of 0.02×0.02 is appropriate for this level of testing.

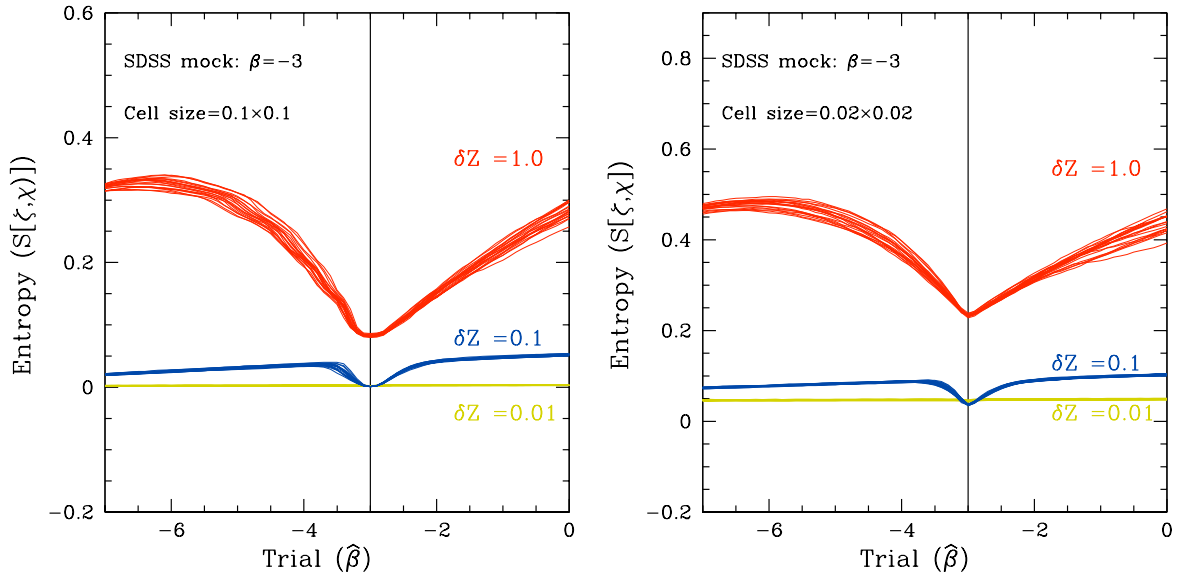


Figure 9.24: Relative entropy, S versus trial $\hat{\beta}$ using SDSS mocks for three different values of δZ : 0.01, 0.1 and 1.0 and varying cell sizes. In each case we have drawn from mocks with a $\beta_{true} = -3.0$. The left-hand panel shows the results for cell size of 0.1×0.1 (100 cells), whereas the right-hand panel is for a cell size 0.02×0.02 (2500 cells). In both cases we observe the that a $\delta Z \gtrsim 0.1$ would be require to recover an adequate signal.

9.3.3 SDSS-Mocks (JTH)

We now test the entropy approach for mocks where the data is very well defined by both a faint and bright limit. To do this we apply our SDDS mocks. Just like MGC, we have explored various scenarios in Figure 9.24 with variations on both the mesh size for the entropy measurement and δZ for ζ to determine the optimum for both. In the left panel of the figure we adopt at cell size of 0.1×0.1 for the mesh and vary δZ for 0.01, 0.1 and 1.0. As with MGC, for $\delta Z = 0.01$ there is no real observable minimisation in the entropy, however, as we increase to $\delta Z \gtrsim 0.1$ we see S clearly minimise for the correction value of $\hat{\beta} = 3.0$ in this case. We observe the same pattern in the right hand pane of Figure 9.24 where the cell size in this case is 0.02×0.02 . For this mesh size there is a clear sharper minimisation compared that of the left panel. Therefore, as we go on to test other values for β_{true} we adopt a $\delta Z = 0.2$ and a cell size of 0.02×0.02 . In Figure 9.25 we have once again apply the same β_{true} values as before: the red lines trace the results for mock with an initial $\beta_{true} = 1.0$, the blue lines are for $\beta_{true} = 2.0$ and the green lines area $\beta_{true} = 3.0$. It is clear that in all three cases the relative entropy once again constrains β_{true} extremely well. Moreover, there does not appear to be any anomalous behaviour that was observed for the same set of mocks under ρ

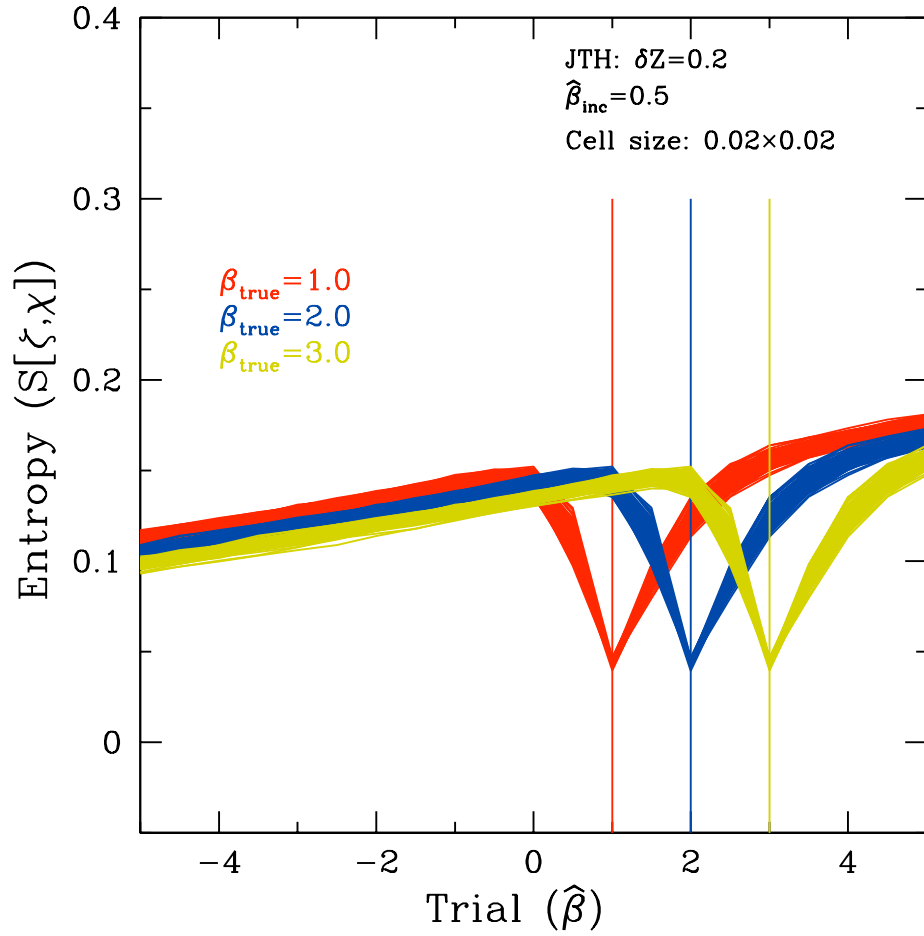


Figure 9.25: Relative entropy, S versus trial $\hat{\beta}$ using SDSS mocks for three different values of β_{true} . For each β_{true} value we have generated 50 mock realisations. In each case we observe the relative entropy constraining β_{true} extremely well.

test.

9.3.3.1 Error Analysis

We apply the same procedure to SDSS as with MGC. However, as we have seen in Figure 9.25, S minimises very sharply for each β_{true} trial, with almost no visual error distribution. We therefore use a very small increment in $\hat{\beta}_{\text{inc}} = 0.005$ in order to accurately assess the error distribution as shown in Figure 9.26. Compared with MGC, we see an even narrower distribution overall with each average $\bar{\beta}$ centred very close to the β_{true} value to within ~ 0.1 . To summarise Figure 9.26 we have for: $\beta_{\text{true}} = 0$ we recover a mean value of $\bar{\beta} = -0.00095$ with a 68% CI of $[-0.015, 0.015]$; $\beta_{\text{true}} = 1$, $\bar{\beta} = 0.99857$ with a 68% CI of $[0.98, 1.015]$; $\beta_{\text{true}} = 2$, $\bar{\beta} = 1.99822$ with a 68% CI of $[1.985, 2.015]$; and $\beta_{\text{true}} = 3$, $\bar{\beta} = 2.9987$ with a 68% CI of $[2.985, 3.015]$.

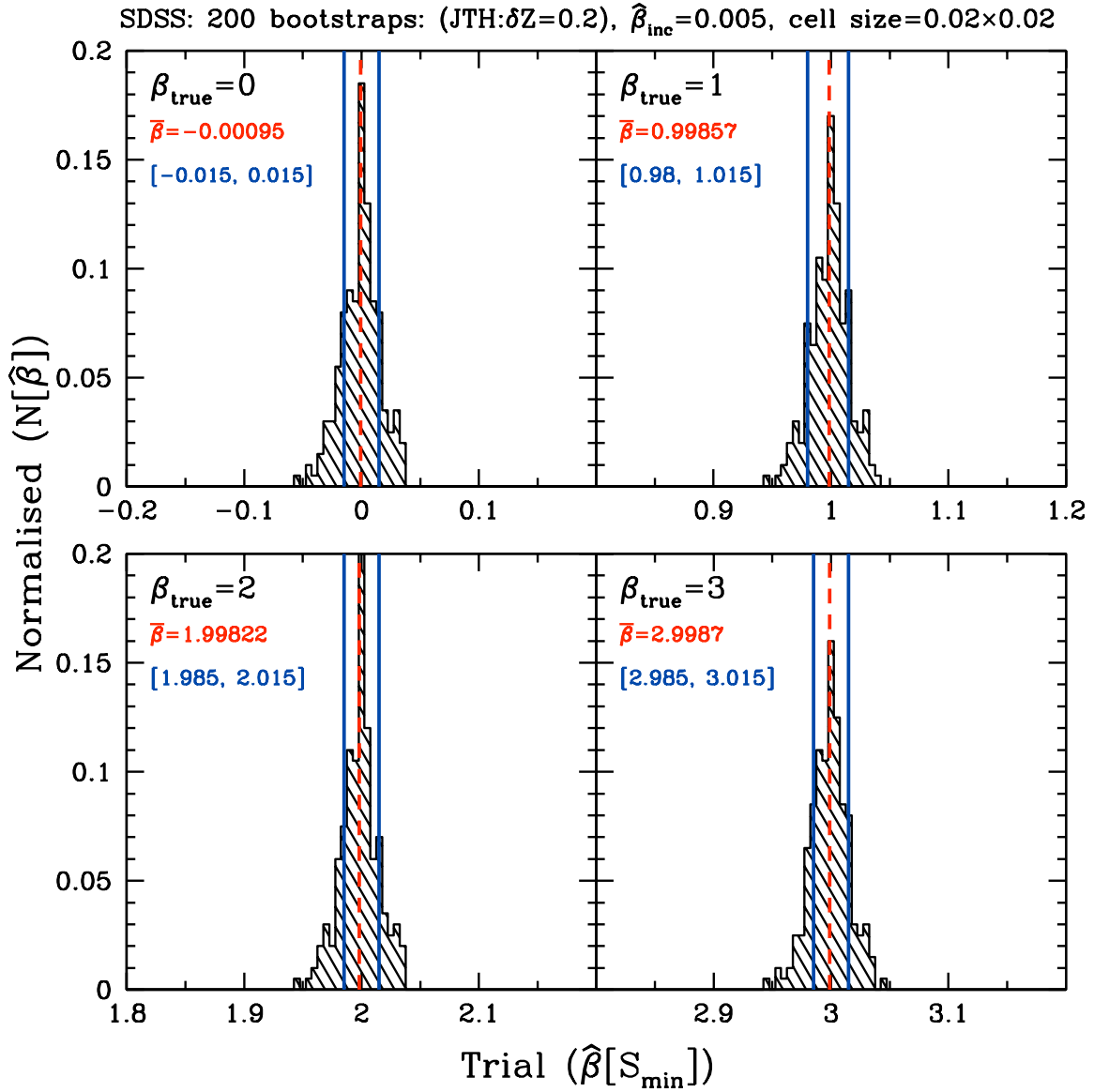


Figure 9.26: $\hat{\beta}(S_{\min})$ distribution for 200 SDSS bootstraps for $\beta_{\text{true}} = 0, 1.0, 2.0$ and 3.0 . The 68% CI's are represented as solid blue vertical lines, and the average of the distribution as a vertical red dashed line.

9.4 Conclusions

In this chapter we have proposed a proof of concept for a new approach to constraining evolution in galaxy redshift surveys. Through the generation of Monte Carlo simulations from the redshift distributions of real galaxy surveys with magnitudes sampled from a β dependent evolutionary model, we have demonstrated that we can effectively constrain evolution. This was achieved by extending our robust statistics to consider

the distribution of ζ and a new random variable, χ . The key to the success can be summarised as follows:

1. The random variable, ζ , has the property of being uniformly distributed on the interval $[0,1]$ for a magnitude-*complete* sample.
2. Our method for generating mock galaxy catalogues with evolution introduces a curve into the apparent magnitude limit that is curved by an amount dependent on the evolutionary corrections to the absolute magnitudes.
3. The construction of our statistics requires us to estimate ζ and χ from an M - Z sample within the apparent magnitude limit defined by a trial magnitude limit m_* . For the purposes of this analysis we have fixed m_* to be straight, and defined it as the faintest galaxy in each corrected data-set. This makes our estimator very sensitive to any curved distribution (caused by evolutionary effects or other) close to the limit.
4. Consequently, correcting the magnitudes of an evolved mock catalogued with the correct value of $\hat{\beta}$ will render the M - Z distribution separable and will be defined with a straight apparent magnitude limit.

To test for correlations we applied two different approaches. The first was the correlation coefficient, ρ and the second, a relative entropy, S , approach. For this analysis we applied our mock catalogues of MGC and SDSS based on the procedure detailed in Chapter 8. We demonstrate that the ρ estimator constrained evolution very well when applied to the MGC mocks. However, when we applied this estimator to the SDSS mocks we found that ρ indicated two distinct ranges for the β parameter. Assuming there are no computational errors, we plan to explore the higher order moments of ρ to see if this can shed any light on this strange result.

However, the relative entropy approach proved to be very successful at constraining β in both the MGC and SDSS cases with very narrow error distributions. Clearly there is scope to develop this approach further to be used on real data. However, this will require a modification of the way in which we sample our mock catalogues. The scenario we have tested in this chapter relies on an M - Z distribution that is representative of an observed survey with realistic evolutionary effects inherent. Our magnitude selection process therefore renders this distribution with a curved apparent magnitude limit. This curving of the limit is not observed in actual survey data. Instead the survey data has an *imposed* straight magnitude cut. Since our procedure exploits the curved

nature of M - Z distribution close to the limit, in its current construction, we do not expect it to recover the corrected evolutionary parameters required to correct the raw magnitudes of a given survey.

The solution to this may lie in the way we model m_* , as discussed in Appendix A, or in the way in which we correct the data before estimating ζ . Since, for a real survey the resulting M - Z distribution is defined by straight limit, we can, as in Rauzy (2001), correct *both* the magnitudes and the distance modulus, Z . This introduces an extra degree of freedom for the way in which the galaxies on M - Z plane can move. Moreover, for any evolutionary corrections that we apply to the data always render the M - Z with a straight limit.

Chapter 10

Discussion and Future Development

“I could tell you what’s happening, but I don’t know if it would really tell you what’s happening.”

Snow, From the movie *Solaris* (2002 version)

10.1 Initial Development of the Completeness Test

Historically, it has often been the case that when estimating luminosity functions, either by parametric or non-parametric means, one assumes,

1. the survey sample data is already complete in apparent magnitude and,
2. the probability distributions $\phi(M)$ and $\rho(z)$ are separable and therefore statistically independent.

Although work by [Efron and Petrosian \(1992, 1999\)](#) provided statistical tools to test the latter assumption of separability, it was [Rauzy \(2001\)](#) (R01) who developed the methodology for a completeness test for magnitude-redshift samples.

10.1.1 Reviving the Rauzy completeness test

We initially revived the work by Rauzy and applied the R01 T_c statistic for the first time to the recent MGC survey sample. We confirmed that the survey data, when not corrected for k - and/or e -corrections, was complete up to the published apparent magnitude limit $m_{\text{lim}} = 20.0$ mag and showed no signs of any residual systematics. With

the addition of $(k+e)$ -corrections we noted that although the T_c statistic confirmed the initial completeness result, there was noticeable (although not statistically significant) systematic drop in the statistic just prior the magnitude limit of the survey. It was concluded that this feature was due to the faint magnitude limit being ‘fuzzied’ by the addition of these corrections.

An interesting point worth noting regarding MGC is that, although there was a published bright limit, the distribution of galaxies on the M - Z plane coupled with our completeness test results indicated that the survey data was well described by a faint limit *only*. This observation would turn out to be crucial in discriminating between surveys that have both a faint limit *and* bright limit that may be more difficult to detect.

10.1.2 The 2dF survey and double truncation

When we applied R01 to the 2dFGRS data, we discovered anomalous behaviour in T_c that initially suggested the survey sample was incomplete. When the data was split according to their APM plates and re-tested, it was found that, for the most part, T_c showed completeness up to the corresponding faint magnitude limits in each plate. It was finally concluded that the total survey data could only be considered complete if secondary *bright* apparent magnitude limit was included and accounted for correctly. Applying R01 to an SDSS-Early Types sample, where the data was published with distinct faint and bright limiting magnitudes, we found similar anomalous behaviour to that of the 2dFGRS. Both of these results lead to our extension of the R01 T_c statistic as published in Johnston, Teodoro & Hendry (2007) - JTH.

The JTH generalisation required the introduction of a quantity, δZ , which fixed the width of the regions S_1 and S_2 to ensure the separable rectangular region could be uniquely defined within the bright and faint limits of a given survey. One would then choose a *suitable* value of δZ that was large enough to overcome effects of shot-noise (see § 10.4), but also small enough to ensure S_1 and S_2 remained within the magnitude limits. By applying the JTH method in this way we demonstrated that both the 2dFGRS and the SDSS-Early Types could now be considered complete up to their respective published faint magnitude limits, $m_{\text{lim}} = 19.45$ mag and $m_{\text{lim}} = 17.45$ mag. Since there was no published bright limit for the 2dFGRS, we adopted a bright limit equal to the brightest galaxy in the subsample, which proved to be adequate.

10.1.3 From T_c to T_v

The T_v statistic, emerged by considering the cumulative distance distribution of the galaxies in a magnitude-redshift survey. The construction of T_v was such that it retained the analogous properties to that of T_c - independence of the spatial distribution of galaxies within the survey and normally distributed with mean zero and variance of unity. Moreover, T_v essentially represents a differential version of the popular V/V_{\max} test. The natural next step developed T_v for surveys with two limiting magnitudes. Therefore, the introduction of the quantity, δM , allowed the separable regions S_3 and S_4 to remain uniquely defined within the two magnitude limits in exactly the same way that δZ was designed for T_c . For comparison, T_v was applied to the three previous surveys - MGC, 2dFGRS and SDSS-Early Types. The results were almost identical to that of the T_c statistic giving further credibility to the reported survey sample selection of all three data-sets.

10.2 Future Work: Part I - Error Estimation

Whilst the JTH method marked a significant improvement over R01, results from 2dFGRS (and the later examination of the CCLQG survey sample) has revealed a possible shortcoming that could potentially undermine the essence of the Rauzy completeness test if not accounted for correctly. As we have already demonstrated throughout this thesis, the completeness test provides an independent approach to help assess and validate the completeness level of a magnitude-redshift sample up to its faint apparent magnitude limit. In the absence of a bright limit, the R01 approach allows the regions S_1 and S_2 , for T_c and S_3 and S_4 , for T_v to grow to a size where the maximum number of galaxies can be included in the respective ζ and τ calculations. This renders both estimators extremely sensitive to any adverse effects inherent in the data, especially close the the faint limit where bias caused by selection effects are more prominent.

In the JTH method we have already discussed that if the choice of δZ and δM are made too small, then the estimators will be dominated by shot noise resulting in a flat-lining of the statistics within the $|3\sigma|$ limits as m_* moves to fainter magnitudes beyond the limit of the survey. The 2dFGRS and CCLQG samples differed to MGC and SDSS in one crucial way - they both have faint magnitude limits that are not well defined by a sharp cut-off. This has the effect whereby, for each incremental increase δZ and δM that overcomes the effects of shot-noise and allows the test statistics to drop below -3σ , we find a range of possible values that indicate the *true* faint magnitude

limit. This in effect would allow any one applying the method to simply choose the magnitude limit to suit their requirements. Therefore, we recognise that to make our statistical tools more robust we require to develop an error estimate for each $T_c(m_*)$ and $T_v(m_*)$ point that will be correlated with the adaptive width (see § 10.4) of the respective δZ and δM .

10.3 The T_v Anomaly

When we were presented with a GALEX selected sample from the Clowes-Campusano Large Quasar Group Survey (CCLQGS) it allowed us to test a data-set that was both deep in its redshift range and utilised photometric redshifts. Once again, it was the application of our tools to *real* data that prompted the need to amend the method. The CCLQG sample provided us with an interesting effect that was only observed in the T_v statistic. This particular sample was described out to a limiting magnitude of $m_{\text{lim}} = 25.8$ mag, however when we applied T_c and T_v there was a noticeable departure in T_v from T_c that lead to an observed peak at $T_v[m_*(22.5)] \approx 18.5\sigma$. Initially, it was thought this may be a result of evolution since this was the first deep survey sample ($z_{\text{max}} \sim 2.5$) to which our statistics had been applied. We had previously discussed the possibility of discrepancies between T_c and T_v where pure luminosity evolution will be a more dominant effect with deeper redshift surveys. Moreover, there were no k - or evolutionary corrections provided with this data-set which could account for the observed effects. However, in this particular case, the differences between T_c and T_v was, in the end, not a result of evolution but instead the nature of the photometric redshifts themselves.

It is well understood among the observational cosmological community that while the use of photometric redshifts over the spectroscopic counterpart is far more efficient for processing large surveys, they severely lack the same level of precision leading to the data being highly rounded. As a result, this creates large discrete artificial gaps in the distribution function for the distance modulus, which is more noticeable for the nearby galaxies. Consequently, this introduced an artificial bias in the T_v estimator which manifested itself as a sharp rise in the resulting T_v curve. We noted that within the statistical literature effects of this nature can be overcome by adding a small amount of uniform random noise to the data. This has the effect of breaking any statistical ties and imposing rank to the data where none was present before. We also noted that this type of effect was also alluded to in R01 when Rauzy considered similar truncation in the magnitude distribution data.

10.4 Future Work: Part II - Optimisation

In Chapter 7 we outlined our proposed method for optimising the T_c and T_v estimators. We detailed an approach that will estimate ζ and τ by assessing their respective signal-to-noise (s/n) levels, $(\zeta/\delta\zeta)$ and $(\tau/\delta\tau)$. We discussed that one of the most efficient ways to incorporate this form of optimisation is to base the ζ and τ calculations on a minimum s/n threshold. This would allow the widths of both δZ and δM to vary for each estimation of ζ and τ - since the number of galaxies increase with increasing distance moduli, the corresponding sampled regions for ζ and τ should, in turn, decrease in size to achieve the required minimum s/n level.

So far, we have generated s/n maps for each of the surveys, and demonstrated how the s/n varies as a function of the trial apparent magnitude limit, m_* and δZ (or δM). One of the main striking features to come out of these maps was a triangular region where δZ or δM has grown too large, rendering the respective regions, S_1 and S_3 too narrow to sample any galaxies within the two defined apparent magnitude limits, m_{lim}^f and m_{lim}^b . Therefore, the s/n was could not be defined, effectively making this a forbidden region. This region in itself is a useful guide to place a maximum limit on either δZ or δM if we are to try and maximise the sampling of the data to test completeness out to as bright a trial magnitude limit as possible. The rest of the maps indicated how the s/n increased as overall as we move out to fainter values of m_* . Although we superimposed white lines showing where the maximum s/n value occurred for each m_* at the corresponding *smallest* δZ and δM widths, in practice, implementing these maps will require a different approach.

Our main goal is to not only optimise the JTH method but also that of the programming code itself in order that the next generation of surveys may be processed more efficiently. This will require us to be able to minimise the width size of the regions that define ζ and τ as much as possible, and still achieve the required target s/n threshold. As an initial test we can adopt a 2 step procedure that represents a simplistic approach to implement these maps into our estimators:

1. For a trial minimum s/n threshold, interpolate across the $(m_*, \delta Z)$ and $(m_*, \delta M)$ range. This will give an average value for the δZ and δM widths at each m_* .
2. Incorporate the interpolation into the T_c and T_v calculation such that both estimators should maintain an approximate constant s/n .

Of course this rudimentary test will be prone to exclude galaxies close the bright magnitude limit of a given survey since for each m_* the width of sampling regions will be

fixed. Further stages of development will allow the code to be more adaptive. However, a balance will need to be struck between how many galaxies we can realistically sample in our calculations, to that of the efficiency of the code and assuring that we are sampling enough to determine if the completeness test is working accurately.

10.5 Using Completeness to Probe Evolution

10.5.1 The mocks

Developing a statistical tool to probe evolutionary models required us to generate a simple but effective methodology for creating Monte Carlo mock galaxy catalogues. For ease of comparison and computer efficiency our mocks were drawn using the observed redshift distributions from MGC, SDSS, and 2dFGRS catalogues. This had the advantage of replicating the observed redshift distribution exactly with in built-clustering and selection effects. We then only had to mock the magnitude distribution which was achieved by sampling absolute magnitudes from known luminosity functions based published luminosity function parameters. This was found to be an effective way to simulate these surveys.

A total of a thousand realisations were generated for each survey, representing both the observer's frame (uncorrected absolute magnitudes) and the galaxy's frame (corrected absolute magnitudes). In terms of the quantities we have been concerned with, both frames of reference are easily summarised as:

$$\text{Galaxy's Frame} \begin{cases} m_{\text{corr}} = m_{\text{obs}} - k(z) - e(z) - A_v(l, b), \\ M_{\text{corr}} = m_{\text{corr}} - 5 \log_{10}(d_L) - 25, \\ Z = 5 \log_{10}(d_L) + 25. \end{cases}$$

$$\text{Observer's Frame} \begin{cases} m = m_{\text{obs}}, \\ M = m - 5 \log_{10}(d_L) - 25, \\ Z = 5 \log_{10}(d_L) + 25, \end{cases}$$

The statistics T_c and T_v were then applied to both frames using the R01 or JTH method where appropriate. Within the observer's frame we found that, in general, both T_c and T_v indicated acceptable levels of completeness up to the survey magnitude limits, m_{lim}^f . With the 2dFGRS and MGC mocks there was a definite rise in T_c and T_v towards the limit for some of the realisations which we concluded may be a sign of mild k - and evolutionary effects arising from the variation of possible magnitude distributions from different sampling seeds.

When examining completeness in the galaxy’s frame we found that both statistics showed completeness for all three surveys entirely consistent with a curved magnitude limit due to the addition of evolutionary corrections.

10.5.2 Testing for evolution

In Chapter 9 we implemented our mock catalogues to probe standard pure luminosity evolution models (PLE). Central to our methodology was the use of the $(\zeta-\chi)$ distribution, where χ represented a new random variable defined directly from the CDF of the redshift distribution. We exploited the fact that ζ and χ are independent, and the resulting joint distribution for a given trial apparent magnitude limit, m_* , brighter than the true limit should be uniformly distributed on a unit square.

Therefore, we utilised our mocks and starting from a Universal LF we introduced and evolved M^* modelled on PLE with a known intrinsic value of β called, β_{true} , and proceeded to sample the mock magnitudes from this evolving LF. Once the mock catalogue was generated we demonstrated that the resulting $(\zeta-\chi)$ distribution indicates correlation. Our goal was then to correct the mock magnitudes with the known PLE model for a series of trial values of $\hat{\beta}_k$. For each corrected data-set we then applied two different statistical techniques to assess the observed correlation between ζ and χ . The first was the correlation coefficient, ρ , whilst the second measured the relative entropy, S , of the distribution. For the correct value of $\hat{\beta}_k$ that equals β_{true} and thus corrects the evolved LF distribution to that of a Universal LF, should result in a $\rho = 0$ for the correlation coefficient and a minimisation of the relative entropy, S .

We applied both approaches to our set of MGC and SDSS Monte Carlo mock catalogues with varying success. In terms of the correlation coefficient, ρ , we found this test statistic to constrain β_{true} fairly well when applied to MGC. We generated 200 bootstrap realisations and found that whilst, for each trial value, $\hat{\beta}$, the overall distribution was roughly centred around the true values for $\rho = 0$, there was a definite skewness in each case. We concluded that this could be attributed to the way in which the M - Z distribution changes its overall curvature near the apparent magnitude limits when we move from negative values to positive values of $\hat{\beta}$.

The results from the SDSS samples yielded surprising results that remain unexplained. Our results showed the following: for each of the 200 bootstraps, ρ crossed zero at two distinct points either side of the β_{true} value. Moreover, ρ appeared to peak sharply significantly above $\rho = 0$ at exactly the β value we were trying to constrain. Assuming there are no programming issues, the reasons for this odd behaviour remains

a baffling one and most definitely requires further exploration. The next step will be to investigate the higher order moments of ρ to see if we can recover any more information.

The application of the entropy approach was a resounding success for constraining β_{true} . Perhaps one of the main advantages this technique has over the ρ is that all the higher order moments are contained within the entropy, S , statistic. For every β_{true} case in both MGC and SDSS, we found that the relative entropy technique minimised at the correct value with no ambiguity at all in the overall $\hat{\beta}$ distribution. In fact, the error distribution about the minimum values of $\hat{\beta}(S_{min})$ were very narrow indeed compared to the results from the ρ , showing very little, if any, skewness.

10.6 Future Work: Part III - Evolution and Other Avenues

Clearly there is still much work to be done in our studies of evolution. For the moment, we have, with the maximum entropy method, a proof of concept that has shown itself to work remarkably well at constraining a single evolutionary parameter. However, further work will adopt this directly to real data where, for uncorrected data, the M - Z distribution is rendered with constant fixed magnitude limits - not curved as we have been using in our tests. This will require a reworking of how we generate our mock catalogues and perhaps, ultimately, draw our mocks from N-body simulations instead of using the observed redshift distributions. Naturally, it would be advantageous to develop this test to include density evolution as well. Ultimately, for real data k and e -corrections are often convolved. Therefore, applying our method to constrain evolution in future surveys will lead a general expression that will convolve PLE and PDE models with k -corrections into a single expression in the form of e.g. a Taylor expansion.

Development of completeness is not restricted to just the magnitude-redshift datasets. There is most certainly scope to adapt and extend all our techniques to other scenarios. Such areas where ROBUST could be extended are bivariate distributions that would include surface brightness and/or colour. Of course, wherever there is a case of testing the independence of two quantities, our statistical analysis could be adapted to suit.

It may also be possible to reverse the role of completeness and incorporate our methodology into the next generation of redshift surveys i.e. we can ask, will it be possible to have completeness by design? This could be an important question as telescope time comes at a premium for the smaller endeavours. If completeness could be incorporated into the observing runs it could potentially maximise the time required

to gather complete samples.

10.7 Concluding Remarks

Statistical cosmology is sure to play an increasingly crucial role for observational cosmology in the not so distant future. For example with the the advent of the Australian SKA Pathfinder and the equivalent MerrKAT in South Africa both acting as the progenitors to the Square Kilometre Array, the cosmology community will surely be overwhelmed with an unprecedented amount of data. Such constructions will mark the next generation in surveys which will be on increasingly adventurous scales (both in volume and sheer numbers of objects). As a result, the need for more sophisticated statistics is quickly becoming a challenge of a dual nature that should not be overlooked - to develop unbiased estimators that can efficiently process large numerical data-sets.

Appendix A

Modelling m_*

If one has an analytical expression for both the $k(z)$ - and evolutionary corrections, should the trial apparent magnitude, m_* , in T_c and T_v be modelled accordingly? Below we consider the alternative ways the M - Z distribution may be constructed and how this impacts our test statistics.

A.1 Methods for Corrections

Approach #1 (Rauzy, 2001): The apparent magnitude for each galaxy was firstly corrected for $k(z)$ -correction, $e(z)$ -correction and extinction, $A(l, b)$. For simplicity, we only consider the effect of combinations of $e(z)$, and $k(z)$ -corrections. Therefore, the corrected apparent magnitude for the i^{th} galaxy, m_i^{corr} is given by,

$$m_i^{corr} = m_i - k(z_i) - e(z_i), \quad (\text{A.1})$$

where m_i is the uncorrected raw apparent magnitude. The author then converts to absolute magnitudes by folding in Equation A.1 to give,

$$\begin{aligned} M_i^{corr} &= m_i^{corr} - 5 \log(d_{L_i}) - 25 \\ &= m_i - 5 \log(d_{L_i}) - 25 - k(z_i) - e(z_i), \end{aligned} \quad (\text{A.2})$$

where, d_{L_i} is the luminosity distance of the i^{th} galaxy. Finally, the same corrections are added to the distance modulus calculations giving,

$$Z_i^{corr} = 5 \log(d_{L_i}) + 25 + k(z_i) + e(z_i) \quad (\text{A.3})$$

Since the corrections are effectively added twice, the corresponding trial limiting apparent magnitude, m_* is modelled as,

$$m_* = M_i^{corr} + Z_i^{corr} = m_*^j \quad (\text{A.4})$$

where m_*^j denotes the initial trial value of m_*^j . This implies that each trial m_*^j will be a fixed straight line on the resulting M - Z plane. As such all the points in this plane will be perturbed in a diagonal direction with the apparent magnitude limits, m_{lim}^b and m_{lim}^f of the survey remaining straight and thus parallel to each m_*^j . This is illustrated in the top panel of Figure A.1 where we show the resulting MGC M - Z distributions where, no corrections have been added (green), only $k(z)$ -corrections have been added (shown in blue), only $e(z)$ -corrections have been added (shown in black), and where both $k(z)$ - and $e(z)$ -corrections have been added (shown in black).

We observe that for the case where the evolutionary correction has only been accounted for, the distribution of points is systematically shifted diagonally downwards compared to that of the green uncorrected data. The effect is small, however, since the value of β in the evolutionary model is small (i.e. $\beta = 0.75$). Moreover, the reason we observe a downward shift is due to the functional form of the model which is given by,

$$e(z) = -\beta \times 2.5 \log_{10}(1 + z_i). \quad (\text{A.5})$$

Where this equation is taken from our analysis of MGC in 3.2.3 (on page 60). In contrast, the addition of the $k(z)$ -corrections, in blue, shows a significant upward shift particularly with the more distant galaxies. When we include both $k(z)$ - and $e(z)$ -corrections, in red, we can see that overall, the distribution is dominated by the $k(z)$ -correction which is not unsurprising.

The movement of galaxies on this M_{corr} - Z_{corr} raises a few very pertinent questions:

1. How sensitive are T_c and T_v in the M_{corr} - Z_{corr} plane?
2. Are the surveys, MGC, SDSS and 2dFGRS simply too shallow for our statistics, in their current construction, to observe effects such as evolution that should, by definition, render the *uncorrected* M - Z distribution un-separable?
3. Consequently, is there an alternative way to model the data that would maximise the sensitivity of our test statistics?

There is no straight forward answer to any of these questions, however, we can begin by addressing the last question first and consider two other possible approaches to modelling m_* .

Approach #2 (Fixing Z): We now consider applying the corrections to the magnitudes *only* and maintain a fixed distance modulus Z that is derived from the redshift distribution that is *not* corrected. With this approach we derive the following expressions for M and Z such that,

$$M_i^{\text{corr}} = m_i - 5 \log(d_{L_i}) - 25 - k(z_i) - e(z_i), \quad (\text{A.6})$$

which is identical to Equation A.2. Since the calculated observed distance modulus, Z_i^{obs} , remains uncorrected throughout the expression is simply given by,

$$Z_i^{\text{obs}} = 5 \log(d_{L_i}) + 25, \quad (\text{A.7})$$

which now leads to trial apparent magnitude limit that could legitimately be modelled as,

$$m_* = M_i^{\text{corr}} + Z_i^{\text{obs}} = m_*^j - k(z_i) - e(z_i). \quad (\text{A.8})$$

If we look at the bottom-left panel of Figure A.1 we can see the resulting M - Z distributions after applying these corrections. By correcting the absolute magnitudes only, we restrict the movement of galaxies on the plane to a horizontal direction. Comparing this to the Approach # 1 where we correct both M and Z , we now observe that if we apply any combination of the $k(z)$ - and $e(z)$ -corrections the resulting distributions shows a distinct curved apparent magnitude limit, m_{lim}^f (depending on the sign convention of the correction). For the case of the $k(z)$ -corrections (in blue), the overall distribution is corrected towards brighter magnitudes compared to that of the case where no corrections have been applied (in green). Conversely, the $e(z)$ -corrections curve the distribution towards fainter magnitudes compared to the uncorrected data. In each case, we observe increased divergence from the uncorrected M - Z distribution as we move to higher values of distance modulus. This is expected since both $k(z)$ - and $e(z)$ -corrections have a strong redshift, z , dependency.

Approach #3 (Fixing M): Conversely we could correct the distance modulus, Z , and maintain a constant uncorrected absolute magnitude. For this case we now have expressions for M and Z given by,

$$M_i^{\text{obs}} = m_i - 5 \log(d_L) - 25. \quad (\text{A.9})$$

Finally, the $(k+e)$ -corrections are added to the distance modulus calculations giving,

$$Z_i^{\text{corr}} = 5 \log(d_L) + 25 + k(z_i) + e(z_i). \quad (\text{A.10})$$

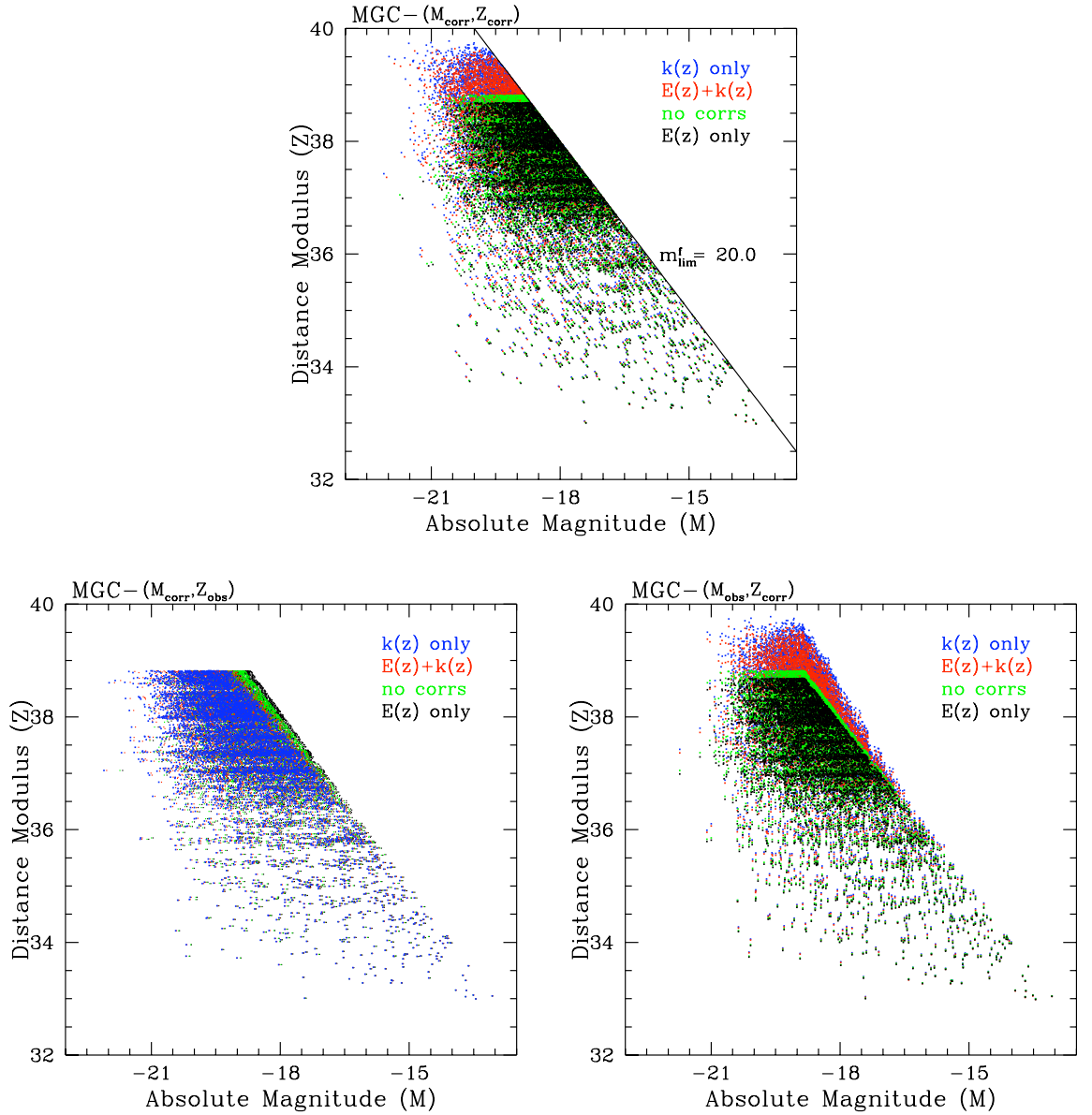


Figure A.1: MGC survey M - Z distributions showing varying methods for correcting the data. In all panels, the distribution of blue points represent $k(z)$ -corrected data, the black shows $e(z)$ -corrected data only, the red shows $(k+e)$ -corrections, and the green points show no corrections to the data. The top panel illustrates Approach # 1, where both M and Z have been corrected. The effect of this sees galaxies move in a diagonal motion which renders m_{lim}^f to remain straight and fixed. The the bottom-left panel shows illustrates Approach # 2 where we the Z distributions remain fixed and we correct only the absolute magnitudes. Galaxies in this cases can move only left or right in the M - Z plane according to the correction model applied to the data. Finally, the bottom panel represents Approach # 3 where the M distribution remains fixed with the Z values being corrected. In this case the corrected galaxies can move only up or down on the plane.

This leads to a trial apparent magnitude limit, m_* that could be modelled as,

$$m_* = M + Z = m_*^j + k(z_i) + e(z_i). \quad (\text{A.11})$$

This form of correction is now shown in the bottom-right panel Figure A.1. In this case galaxies can only move vertically on the M - Z plane creating, once again, a curving of the magnitude limit.

A.2 How Should We Sample the Data?

A.2.1 A constant m_*

Now that we have established a further two approaches for correcting redshift-magnitude survey data, we can now examine their impact on the T_c and T_v completeness statistics. Although the three possible routes one can choose to obtain an M - Z distribution that represents a corrected data-set, we can still construct our test statistics based on an m_* that remains modelled as a straight line in all three cases. In Figure A.2 we apply m_* in exactly this way to all three approaches detail in the previous section.

In the top-left panel we show the T_c and T_v results from Chapter 3 (page 62) where no corrections have been added (i.e. the observed raw data). In the panel right of this we show results where corrections have been added to both M and Z as in the original R01 procedure (Approach #1). You will see the red lines in this panel representing the same $(k+e)$ -corrections also discussed in Chapter 3 (page 62). We also show the effect of the individual corrections of evolution and k shown respectively in black and blue. The effect of the e -correction alone seems, to the eye, to make only a slight difference in the M - Z distribution compared to the uncorrected case shown in the top panel of Figure A.1. However, we observe a significant departure from the 3σ completeness level at $m_* \sim 18.1$ mag as indicated in Figure A.2 compared to the k -correction.

If we look now at the bottom-left panel in Figure A.2, where we apply Approach #2, we now observe very different behaviour. By sampling the galaxies within a region defined by an m_* that is straight and not modelled by the k - and/or e -corrections, we obtain a new measure of completeness. It is this approach that was successfully applied in Chapter 9 when creating our mock galaxy catalogues to probe evolutionary models. Since the corrections cause the magnitude limit of the M - Z distribution to curve according the particular model applied, we observe T_c and T_v to reflect this. In the bottom-left panel of Figure A.2 the evolution corrected magnitudes (in black) recover a T_c and T_v that indicate the data is complete up to the survey limit of $m_{\text{lim}} = 20.0$ mag.

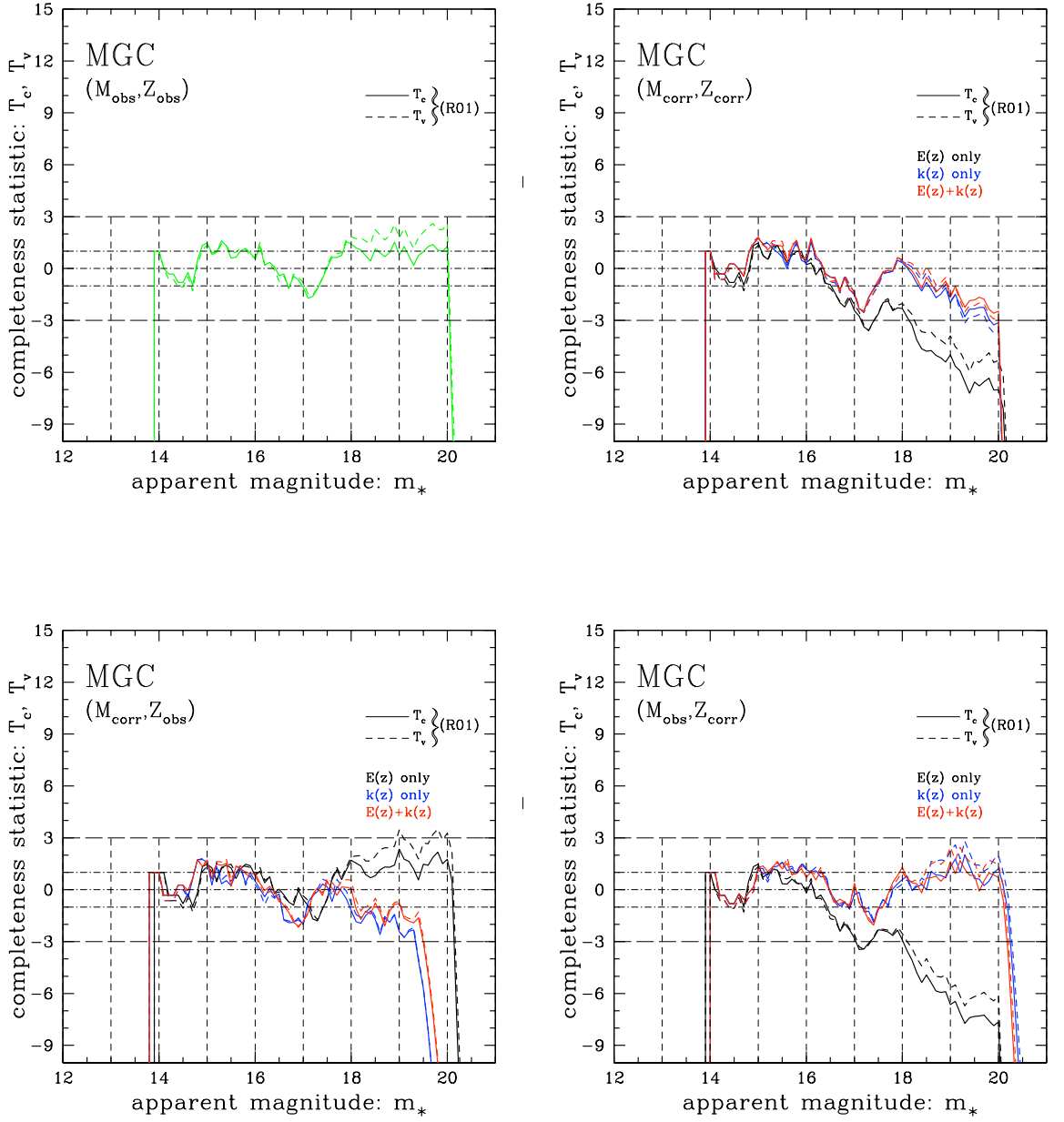


Figure A.2: MGC survey T_c and T_v statistics resulting from applying varying methods for correcting the data. In each case we have assume a faint limit only (i.e R01 method). The top-left panel shows the results where no corrections have been added to either the magnitudes or the distance modulus. In the remaining panels the black lines are resultant from applying $e(z)$ -correction only, the blue represent the $k(z)$ -corrected data and the red lines represent the combined $(k + e)$ -corrections. To distinguish between T_c and T_v we have use a solid line and a dashed line respectively. In the top-right panel, both the distance modulus, Z , and the absolute magnitudes, M , have been corrected. In the bottom-left panel we have corrected M only, and in the bottom-right pane we correct Z only.

However, both T_c and T_v do not drop sharply until $m_* \sim 20.1$ mag. Conversely, when the k and $(k+e)$ -corrections are now included both T_c and T_v drop at respective apparent magnitude limits of $m_* \sim 19.5$ mag (blue lines) and $m_* \sim 19.7$ mag (red lines). As discussed in Chapter 9, this behaviour reflects the direction of the curved M - Z distribution. If, in the case of the e -correction, the M - Z distribution is curved towards fainter magnitudes than the observed uncorrected data, then it is inevitable that T_c and T_v will determine the limit to be approximately equal to faintest galaxy in the corrected data-set. Since the corrections we are considering are z -dependent, the rate at which the distribution is curved becomes more acute with increasing redshift as we have already discussed for Figure A.1. Moreover, as we move to greater redshifts we are estimating T_c and T_v with a substantially increased number of galaxies contained within a larger volume. Therefore, galaxies contained within this region are more likely to dominate the statistical results compare to the contribution of nearby galaxies. As such, when we observe T_c and T_v begin to drop before the actual m_{lim} of the survey as with the k and $(k+e)$ -corrections, the galaxies at higher distance moduli are now significantly brighter than $m_* = m_{\text{lim}}^f = 20.0$ mag. This manifests as a deficit in the number of galaxies counted in the S_2 and S_4 regions for ζ and τ respectively, leading to T_c and T_v dropping systematically below the -3σ limit.

The same behaviour is observed when we now correct the distance moduli and keep the absolute magnitudes as ‘raw’ as in Approach #3 (see the bottom-right panel in Figure A.2). Since the corrections are begin *added* to Z , the behaviour of T_c and T_v is the opposite to that of Approach #2 where the correction are subtracted from the absolute magnitudes.

The fundamental construction of the T_c and T_v statistics do require further investigation if we are to improve and maximise their use within the analysis of redshift surveys (and perhaps other areas of cosmology). Therefore, exploration into modelling m_* , when applying Approaches #2 and 3, may lead us to a better understanding of our statistics. This would be achieved by modelling m_* according to the evolutionary and/or k -correction used to correct the data. Thus, as the magnitude limit in the M - Z distribution curves for a given model correction, the m_* line would also curve according to the z -dependent model.

Bibliography

- Abazajian, K., *et al.* : 2003, *Astronomical Journal* **126**, 2081
- Abazajian, K. e.: 2009, *Astrophysical Journal, Supplement* **182**, 543
- Abraham, R. G., *et al.* : 2004, *Astronomical Journal* **127**, 2455
- Abramowitz, M. and Stegun, I.: 1964, *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, Dover Publications
- Adelman-McCarthy, *et al.*: 2006, *Astrophysical Journal, Supplement* **162**, 38
- Adelman-McCarthy, J. K. and *et al.* : 2007, *Astrophysical Journal, Supplement* **172**, 634
- Aldrich, J.: 1997, *Statistical Science* **12(3)**, 162
- Avni, Y. and Bahcall, J. N.: 1980, *Astrophysical Journal* **235**, 694
- Ball, N. M., Brunner, R. J., Myers, A. D., Strand, N. E., Alberts, S. L., and Tchong, D.: 2008, *Astrophysical Journal* **683**, 12
- Baum, W. A.: 1962, in G. C. McVittie (ed.), *Problems of Extra-Galactic Research*, Vol. 15 of *IAU Symposium*, pp 390–+
- Beckwith, S. V. W., *et al.* : 2006, *Astronomical Journal* **132**, 1729
- Benítez, N.: 2000, *Astrophysical Journal* **536**, 571
- Berlind, A. A. and Weinberg, D. H.: 2002, *Astrophysical Journal* **575**, 587
- Bernardi, M. *et al.* : 2003a, *Astronomical Journal* **125**, 1817
- Bernardi, M. *et al.* : 2003b, *Astronomical Journal* **125**, 1849

- Bernardi, M., Sheth, R. K., Nichol, R. C., Schneider, D. P., and Brinkmann, J.: 2005, *Astronomical Journal* **129**, 61
- Bertin, E. and Arnouts, S.: 1996, *Astronomy and Astrophysics, Supplement* **117**, 393
- Binggeli, B., Sandage, A., and Tammann, G. A.: 1988, *Annual Review of Astronomy and Astrophysics* **26**, 509
- Blake, C., *et al.* : 2008, *Astronomy and Geophysics* **49(5)**, 050000
- Blanton, M. R., Lin, H., Lupton, R. H., Maley, F. M., Young, N., Zehavi, I., and Loveday, J.: 2003, *Astronomical Journal* **125**, 2276
- Bolzonella, M., Miralles, J.-M., and Pelló, R.: 2000, *Astronomy and Astrophysics* **363**, 476
- Boyle, B. J., Fong, R., Shanks, T., and Peterson, B. A.: 1990, *Monthly Notices of the Royal Astronomical Society* **243**, 1
- Brunner, R. J., Connolly, A. J., Szalay, A. S., and Bershad, M. A.: 1997, *Astrophysical Journal, Letters* **482**, L21+
- Budavári, T.: 2009, *Astrophysical Journal* **695**, 747
- Caditz, D. and Petrosian, V.: 1993, *Astrophysical Journal* **416**, 450
- Cannon, R., *et al.* : 2006, *Monthly Notices of the Royal Astronomical Society* **372**, 425
- Carlberg, R. G., *et al.* : 1999, *Royal Society of London Philosophical Transactions Series A* **357**, 167
- Charlier, C. V. L.: 1922, *Arkiv. för. Astron. Fys.*, **16**, 1
- Choloniewski, J.: 1986, *Monthly Notices of the Royal Astronomical Society* **223**, 1
- Choloniewski, J.: 1987, *Monthly Notices of the Royal Astronomical Society* **226**, 273
- Christensen, C. G.: 1975, *Astronomical Journal* **80**, 282
- Cimatti, A., *et al.* : 2002, *Astronomy and Astrophysics* **381**, L68
- Colless, M., *et al.* : 2001, *Monthly Notices of the Royal Astronomical Society* **328**, 1039

- Colless, M.: 1998, in S. Colombi, Y. Mellier, and B. Raban (eds.), *Wide Field Surveys in Cosmology, 14th IAP meeting held May 26-30, 1998, Paris. Publisher: Editions Frontieres. ISBN: 2-8 6332-241-9, p. 77.*, pp 77–+
- Connolly, A. J., Csabai, I., Szalay, A. S., Koo, D. C., Kron, R. G., and Munn, J. A.: 1995, *Astronomical Journal* **110**, 2655
- Croom, S. M., Smith, R. J., Boyle, B. J., Shanks, T., Loaring, N. S., Miller, L., and Lewis, I. J.: 2001, *Monthly Notices of the Royal Astronomical Society* **322**, L29
- Cross, N., e.: 2001, *Monthly Notices of the Royal Astronomical Society* **324**, 825
- Cross, N. J. G., Driver, S. P., Liske, J., Lemon, D. J., Peacock, J. A., Cole, S., Norberg, P., and Sutherland, W. J.: 2004, *Monthly Notices of the Royal Astronomical Society* **349**, 576
- da Costa, L. N., *et al.* : 1988, *Astrophysical Journal* **327**, 544
- da Costa, L. N., *et al.* : 1998, *Astronomical Journal* **116**, 1
- Davis, M. *et al.* : 2003, *Proc. SPIE* **4834**, 161
- Davis, M. and Huchra, J.: 1982, *Astrophysical Journal* **254**, 437
- Davis, M., Tonry, J., Huchra, J., and Latham, D. W.: 1980, *Astrophysical Journal, Letters* **238**, L113
- de Lapparent, V., Geller, M. J., and Huchra, J. P.: 1989, *Astrophysical Journal* **343**, 1
- Driver, S. P., Liske, J., Cross, N. J. G., De Propris, R., and Allen, P. D.: 2005, *Monthly Notices of the Royal Astronomical Society* **360**, 81
- Driver, S. P., Phillipps, S., Davies, J. I., Morgan, I., and Disney, M. J.: 1994, *Monthly Notices of the Royal Astronomical Society* **266**, 155
- Drory, N., *et al.* : 2000, in A. Mazure, O. Le Fèvre, and V. Le Brun (eds.), *Clustering at High Redshift*, Vol. 200 of *Astronomical Society of the Pacific Conference Series*, pp 91–+
- Dunkley, J., *et al.* : 2009, *Astrophysical Journal, Supplement* **180**, 306
- Eales, S.: 1993, *Astrophysical Journal* **404**, 51

- Efron, B. and Petrosian, V.: 1992, *Astrophysical Journal* **399**, 345
- Efron, B. and Petrosian, V.: 1999, *Journal of the American Statistical Association* **94(447)**, 824
- Efstathiou, G., Ellis, R. S., and Peterson, B. A.: 1988, *Monthly Notices of the Royal Astronomical Society* **232**, 431
- Einstein, A.: 1915, *Sitzungsberichte der Königlich Preußischen Akademie der Wissenschaften (Berlin)*, Seite 844-847. pp 844–847
- Eisenstein, D. J., *et al.* : 2001, *Astronomical Journal* **122**, 2267
- Ellis, R. S., Colless, M., Broadhurst, T., Heyl, J., and Glazebrook, K.: 1996, *Monthly Notices of the Royal Astronomical Society* **280**, 235
- Felten, J. E.: 1976, *Astrophysical Journal* **207**, 700
- Felten, J. E.: 1977, *Astronomical Journal* **82**, 861
- Fernández-Soto, A., Lanzetta, K. M., and Yahil, A.: 1999, *Astrophysical Journal* **513**, 34
- Fisher, J. R. and Tully, R. B.: 1981, *Astrophysical Journal, Supplement* **47**, 139
- Fisher, R. A.: 1912, *Messenger of Mathematics* **41**, 155
- Fisher, R. A.: 1922, *Philos. Trans. Roy. Soc. London Ser. A* **222**, 309
- Francis, P. J., Nelson, B. O., and Cutri, R. M.: 2004, *Astronomical Journal* **127**, 646
- Fukugita, M., Ichikawa, T., Gunn, J. E., Doi, M., Shimasaku, K., and Schneider, D. P.: 1996, *Astronomical Journal* **111**, 1748
- Gamow, G.: 1946, *Physical Review* **70**, 572
- Geller, M. J. and Huchra, J. P.: 1989, *Science* **246**, 897
- Gold, B., *et al.* : 2009, *Astrophysical Journal, Supplement* **180**, 265
- Gregory, S. A. and Thompson, L. A.: 1978a, *Nature* **274**, 450
- Gregory, S. A. and Thompson, L. A.: 1978b, *Astrophysical Journal* **222**, 784
- Gregory, S. A. and Thompson, L. A.: 1984, *Astrophysical Journal* **286**, 422

- Gunn, J. E., *et al.* : 1998, *Astronomical Journal* **116**, 3040
- Hao, L. and Strauss, *et al.* : 2005, *Astronomical Journal* **129**, 1795
- Heyl, J., Colless, M., Ellis, R. S., and Broadhurst, T.: 1997, *Monthly Notices of the Royal Astronomical Society* **285**, 613
- Hinshaw, G., *et al.* : 2009, *Astrophysical Journal, Supplement* **180**, 225
- Hogg, D. W., Baldry, I. K., Blanton, M. R., and Eisenstein, D. J.: 2002, *ArXiv*: astro-ph/0210394
- Hogg, D. W., Finkbeiner, D. P., Schlegel, D. J., and Gunn, J. E.: 2001, *Astronomical Journal* **122**, 2129
- Hubble, E.: 1929, *Proceedings of the National Academy of Science* **15**, 168
- Hubble, E.: 1936a, *Astrophysical Journal* **84**, 517
- Hubble, E. P.: 1936b, *Realm of the Nebulae*, Realm of the Nebulae, by E.P. Hubble. New Haven: Yale University Press, 1936
- Huchra, J., Davis, M., Latham, D., and Tonry, J.: 1983, *Astrophysical Journal, Supplement* **52**, 89
- Huchra, J. and Sargent, W. L. W.: 1973, *Astrophysical Journal* **186**, 433
- Hudson, M. J. and Lynden-Bell, D.: 1991, *Monthly Notices of the Royal Astronomical Society* **252**, 219
- Humason, M. L., Mayall, N. U., and Sandage, A. R.: 1956, *Astronomical Journal* **61**, 97
- Impey, C. D., Sprayberry, D., Irwin, M. J., and Bothun, G. D.: 1996, *Astrophysical Journal, Supplement* **105**, 209
- Jackson, J. C.: 1974, *Monthly Notices of the Royal Astronomical Society* **166**, 281
- Johnston, R., Teodoro, L., and Hendry, M.: 2007, *Monthly Notices of the Royal Astronomical Society* **376**, 1757
- Jones, D. H., Saunders, W., Read, M., and Colless, M.: 2005, *Publications of the Astronomical Society of Australia* **22**, 277

- Kafka, P.: 1967, *Nature* **213**, 346
- Kaiser, N., *et al.* : 2002, in J. A. Tyson and S. Wolff (eds.), *Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series*, Vol. 4836 of *Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series*, pp 154–164
- Kelly, B. C., Fan, X., and Vestergaard, M.: 2008, *Astrophysical Journal* **682**, 874
- Kiang, T.: 1961, *Monthly Notices of the Royal Astronomical Society* **122**, 263
- Kirshner, R. P., Oemler, Jr., A., and Schechter, P. L.: 1978, *Astronomical Journal* **83**, 1549
- Kirshner, R. P., Oemler, Jr., A., and Schechter, P. L.: 1979, *Astronomical Journal* **84**, 951
- Kirshner, R. P., Oemler, Jr., A., Schechter, P. L., and Shectman, S. A.: 1981, *Astrophysical Journal, Letters* **248**, L57
- Kirshner, R. P., Oemler, Jr., A., Schechter, P. L., and Shectman, S. A.: 1983, *Astronomical Journal* **88**, 1285
- Kirshner, R. P., Oemler, A. J., Schechter, P. L., and Shectman, S. A.: 1987, *Astrophysical Journal* **314**, 493
- Kodama, T., Bell, E. F., and Bower, R. G.: 1999, *Monthly Notices of the Royal Astronomical Society* **302**, 152
- Komatsu, E., *et al.* : 2009, *Astrophysical Journal, Supplement* **180**, 330
- Koo, D. C.: 1985, *Astronomical Journal* **90**, 418
- Kullback, S. and Leibler, R. A.: 1951, *Annals of Mathematical Statistics* **22**, 79
- Lawrence, A., *et al.* : 2007, *Monthly Notices of the Royal Astronomical Society* **379**, 1599
- Lawrence, A., Walker, D., Rowan-Robinson, M., Leech, K. J., and Penston, M. V.: 1986, *Monthly Notices of the Royal Astronomical Society* **219**, 687
- Le Fèvre, O., *et al.* : 2004, *Astronomy and Astrophysics* **428**, 1043
- Lilly, S. J. *et al.* : 2007, *Astrophysical Journal, Supplement* **172**, 70

- Lilly, S. J., Le Fevre, O., Crampton, D., Hammer, F., and Tresse, L.: 1995, *Astrophysical Journal* **455**, 50
- Lima, M., Cunha, C. E., Oyaizu, H., Frieman, J., Lin, H., and Sheldon, E. S.: 2008, *Monthly Notices of the Royal Astronomical Society* **390**, 118
- Liske, J., Lemon, D. J., Driver, S. P., Cross, N. J. G., and Couch, W. J.: 2003, *Monthly Notices of the Royal Astronomical Society* **344**, 307
- Loh, E. D. and Spillar, E. J.: 1986a, *Astrophysical Journal, Letters* **307**, L1
- Loh, E. D. and Spillar, E. J.: 1986b, *Astrophysical Journal* **303**, 154
- Loveday, J., Peterson, B. A., Efstathiou, G., and Maddox, S. J.: 1992, *Astrophysical Journal* **390**, 338
- Lupton, R., *et al.* : 2001, in F. R. Harnden, Jr., F. A. Primini, and H. E. Payne (eds.), *ASP Conf. Ser. 238: Astronomical Data Analysis Software and Systems X*, pp 269–+
- Lynden-Bell, D.: 1971, *Monthly Notices of the Royal Astronomical Society* **155**, 95
- Maccacaro, T., della Ceca, R., Gioia, I. M., Morris, S. L., Stocke, J. T., and Wolter, A.: 1991, *Astrophysical Journal* **374**, 117
- Maddox, S. J., Efstathiou, G., Sutherland, W. J., and Loveday, J.: 1990, *Monthly Notices of the Royal Astronomical Society* **243**, 692
- Madgwick, D. S., *et al.* : 2002, *Monthly Notices of the Royal Astronomical Society* **333**, 133
- Maloney, A. and Petrosian, V.: 1999, *Astrophysical Journal* **518**, 32
- Markarian, B. E.: 1967, *Astrofizika* **3**, 55
- Markarian, B. E.: 1969a, *Astrofizika* **5**, 443
- Markarian, B. E.: 1969b, *Astrofizika* **5**, 581
- Markarian, B. E., Lipovetskij, V. A., and Lipovetsky, V. A.: 1971, *Astrofizika* **7**, 511
- Marshall, H. L., Tananbaum, H., Avni, Y., and Zamorani, G.: 1983, *Astrophysical Journal* **269**, 35

- Nicoll, J. F. and Segal, I. E.: 1983, *Astronomy and Astrophysics* **118**, 180
- Nolta, M. R., *et al.* : 2009, *Astrophysical Journal, Supplement* **180**, 296
- Norberg, P., *et al.* : 2001, *Monthly Notices of the Royal Astronomical Society* **328**, 64
- Norberg, P., *et al.* : 2002a, *Monthly Notices of the Royal Astronomical Society* **336**, 907
- Norberg, P., *et al.* : 2002b, *Monthly Notices of the Royal Astronomical Society* **332**, 827
- Outram, P. J., Hoyle, F., Shanks, T., Croom, S. M., Boyle, B. J., Miller, L., Smith, R. J., and Myers, A. D.: 2003, *Monthly Notices of the Royal Astronomical Society* **342**, 483
- Oyaizu, H., Lima, M., Cunha, C. E., Lin, H., and Frieman, J.: 2008, *Astrophysical Journal* **689**, 709
- Page, M. J. and Carrera, F. J.: 2000, *Monthly Notices of the Royal Astronomical Society* **311**, 433
- Peacock, J. A. *et al.* : 2001, *Nature* **410**, 169
- Peebles, P. J. E. and Yu, J. T.: 1970, *Astrophysical Journal* **162**, 815
- Percival, W. J. *et al.* : 2001, *Monthly Notices of the Royal Astronomical Society* **327**, 1297
- Perlmutter, S., Turner, M. S., and White, M.: 1999, *Physical Review Letters* **83**, 670
- Petrosian, V.: 1992, in *Statistical Challenges in Modern Astronomy*, pp 173–200
- Petrosian, V.: 2002, in R. F. Green, E. Y. Khachikian, and D. B. Sanders (eds.), *IAU Colloq. 184: AGN Surveys*, Vol. 284 of *Astronomical Society of the Pacific Conference Series*, pp 389–+
- Pier, J. R., Munn, J. A., Hindsley, R. B., Hennessy, G. S., Kent, S. M., Lupton, R. H., and Ivezić, Ž.: 2003, *Astronomical Journal* **125**, 1559
- Poggianti, B. M.: 1997, *Astronomy and Astrophysics, Supplement* **122**, 399
- Press, W. H. and Schechter, P.: 1974, *Astrophysical Journal* **187**, 425

- Qin, Y. P. and Xie, G. Z.: 1997, *Astrophysical Journal* **486**, 100
- Qin, Y.-P. and Xie, G.-Z.: 1999, *Astronomy and Astrophysics* **341**, 693
- Ratcliffe, A., *et al.* : 1996, *Monthly Notices of the Royal Astronomical Society* **281**, L47+
- Rauzy, S.: 2001, *Monthly Notices of the Royal Astronomical Society* **324**, 51
- Rauzy, S., Hendry, M. A., and D'Mellow, K.: 2001, *Monthly Notices of the Royal Astronomical Society* **328**, 1016
- Richards, G. T. *et al.* : 2006, *Astronomical Journal* **131**, 2766
- Riess, A. G., *et al.* : 1998, *Astronomical Journal* **116**, 1009
- Rossi, G. and Sheth, R. K.: 2008, *Monthly Notices of the Royal Astronomical Society* **387**, 735
- Sachs, R. K. and Wolfe, A. M.: 1967, *Astrophysical Journal* **147**, 73
- Salzer, J. J. and Haynes, M. P.: 1996, in E. D. Skillman (ed.), *The Minnesota Lectures on Extragalactic Neutral Hydrogen*, Vol. 106 of *Astronomical Society of the Pacific Conference Series*, pp 357–+
- Sandage, A. and Tammann, G. A.: 1981, *Journal of the British Astronomical Association* **91**, 522
- Sandage, A., Tammann, G. A., and Yahil, A.: 1979, *Astrophysical Journal* **232**, 352
- Santiago, B., Strauss, M., Lahav, O., Davis, M., Dressler, A., and Huchra, J.: 1996, *Astrophysical Journal* **461**, 38
- Saunders, W., *et al.* : 1990, *Monthly Notices of the Royal Astronomical Society* **242**, 318
- Saunders, W., *et al.* : 2000, *Monthly Notices of the Royal Astronomical Society* **317**, 55
- Schafer, C. M.: 2007, *Astrophysical Journal* **661**, 703
- Schechter, P.: 1976, *Astrophysical Journal* **203**, 297

- Schlegel, D. J., Finkbeiner, D. P., and Davis, M.: 1998, *Astrophysical Journal* **500**, 525
- Schmidt, M.: 1968, *Astrophysical Journal* **151**, 393
- Schmidt, M.: 1972, *Astrophysical Journal* **176**, 303
- Schmidt, M.: 1976, *Astrophysical Journal, Letters* **209**, L55+
- Schmitt, J. H. M. M.: 1990, *Astronomy and Astrophysics* **240**, 556
- Schneider, S. E.: 1996, in E. D. Skillman (ed.), *The Minnesota Lectures on Extragalactic Neutral Hydrogen*, Vol. 106 of *Astronomical Society of the Pacific Conference Series*, pp 323–+
- Seljak, U.: 2000, *Monthly Notices of the Royal Astronomical Society* **318**, 203
- Shapiro, S. L.: 1971, *Astronomical Journal* **76**, 291
- Shectman, S. A., Landy, S. D., Oemler, A., Tucker, D. L., Lin, H., Kirshner, R. P., and Schechter, P. L.: 1996, *Astrophysical Journal* **470**, 172
- Sheth, R. K.: 2007, *Monthly Notices of the Royal Astronomical Society* **378**, 709
- Silk, J.: 1967, *Nature* **215**, 1155
- Sivia, D. S. and Skilling, J.: 2006, *Data Analysis - A Bayesian Tutorial*, Oxford Science Publications, Page 183, second edition
- Skrutskie, M. F., *et al.* : 2006, *Astronomical Journal* **131**, 1163
- Slipher, V. M.: 1913, *Lowell Observatory Bulletin* **2**, 56
- Smith, J. A., *et al.* : 2002, *Astronomical Journal* **123**, 2121
- Smoot, G. F., *et al.* : 1992, *Astrophysical Journal, Letters* **396**, L1
- Sodre, L. J. and Lahav, O.: 1993, *Monthly Notices of the Royal Astronomical Society* **260**, 285
- Somerville, R. S., Lee, K., Ferguson, H. C., Gardner, J. P., Moustakas, L. A., and Giavalisco, M.: 2004, *Astrophysical Journal, Letters* **600**, L171
- Spergel, D. N., *et al.* : 2003, *Astrophysical Journal, Supplement* **148**, 175

- Springel, V., *et al.* : 2005, *Nature* **435**, 629
- Springel, V. and White, S. D. M.: 1998, *Monthly Notices of the Royal Astronomical Society* **298**, 143
- Stabenau, H. F., Connolly, A., and Jain, B.: 2008, *Monthly Notices of the Royal Astronomical Society* **387**, 1215
- Steidel, C. C., Adelberger, K. L., Shapley, A. E., Pettini, M., Dickinson, M., and Giavalisco, M.: 2003, *Astrophysical Journal* **592**, 728
- Stoughton, C., *et al.* : 2002, *Astronomical Journal* **123**, 485
- Strauss, M. A., *et al.* : 2002, *Astronomical Journal* **124**, 1810
- Subbarao, M. U., Connolly, A. J., Szalay, A. S., and Koo, D. C.: 1996, *Astronomical Journal* **112**, 929
- Takeuchi, T. T., Yoshikawa, K., and Ishii, T. T.: 2000, *Astrophysical Journal, Supplement* **129**, 1
- Taylor, A. R. and Braun, R. (eds.): 1999, *Science with the Square Kilometer Array : a next generation world radio observatory*.
- Tegmark, M., Hamilton, A. J. S., and Xu, Y.: 2002, *Monthly Notices of the Royal Astronomical Society* **335**, 887
- Turner, E. L.: 1979, *Astrophysical Journal* **231**, 645
- Turner, M. S.: 1998, *APS Meeting Abstracts* pp 101–+
- van den Bergh, S.: 1961, *Zeitschrift fur Astrophysik* **53**, 219
- van Waerbeke, L., Mathez, G., Mellier, Y., Bonnet, H., and Lachize-Rey, M.: 1996, *Astronomy and Astrophysics* **316**, 1
- Vanzella, E., *et al.* : 2005, *Astronomy and Astrophysics* **434**, 53
- Vettolani, G., *et al.* : 1997, *Astronomy and Astrophysics* **325**, 954
- Wang, Y., Bahcall, N., and Turner, E. L.: 1998, *Astronomical Journal* **116**, 2081
- Williams, R. E., *et al.* : 1996, *Astronomical Journal* **112**, 1335

Williams, R. E., *et al.* : 2000, *Astronomical Journal* **120**, 2735

Willmer, C. N. A.: 1997, *Astronomical Journal* **114**, 898

Wittman, D.: 2009, *ArXiv*: 0905.0892v2

Wolf, C.: 2009, *ArXiv*: 0904.3438

Wright, E. L., *et al.* : 1992, *Astrophysical Journal, Letters* **396**, L13

York, D. G., *et al.* : 2000, *Astronomical Journal* **120**, 1579

Zwaan, M. A., Briggs, F. H., Sprayberry, D., and Sorar, E.: 1997, *Astrophysical Journal* **490**, 173